

BiMol-Diff: A Unified Diffusion Framework for Molecular Generation and Captioning

Aditya Hemant Shahane^{1*} Anuj Kumar Sirohi^{1*} Devansh Arora²
Nitin Kumar² Prathosh A P^{3,4} Sandeep Kumar¹

¹Indian Institute of Technology Delhi, ²Indian Institute of Technology Ropar

³Indian Institute of Science Bengaluru, ⁴Latentforce.ai

Abstract

Bridging molecular structures and natural language is essential for controllable design. Autoregressive models struggle with long-range dependencies, while standard diffusion processes apply uniform corruption across positions, which can distort structurally informative tokens. We present BiMol-Diff, a unified diffusion framework for the paired tasks of text-conditioned molecule generation and molecule captioning. Our key component is a *Token-aware noise schedule* that assigns position-dependent corruption based on token recovery difficulty, preserving harder-to-recover substructures during the forward process. On ChEBI-20 and M3-20M, BiMol-Diff improves molecule reconstruction with a 15.4% relative gain in Exact Match and achieves strong captioning results, attaining the best BLEU and BERTScore among compared baselines. These results indicate token-aware noising improves fidelity in molecular structure–language modeling. Code link [GitHub](#).

1 Introduction

Designing molecules requires balancing two distinct modalities: the molecular *structure* of atoms and bonds, and the linguistic descriptions of molecular properties (e.g., “a kinase inhibitor with good solubility”). Bridging this gap through cross-modal translation enables chemically grounded workflows where designs can be iteratively refined via natural language (Zheng et al., 2023; Fatemi et al., 2024). This motivates two complementary directions: (i) **molecule generation** (*text*→*molecule*) to materialize descriptions into structures, and (ii) **molecule captioning** (*molecule*→*text*) to explain molecular structure and properties via natural-language descriptions (Edwards et al., 2022; Zeng et al., 2024). However, despite their symmetry, these tasks are typically developed in isolation with separate mod-

els and objectives, limiting consistency and reuse in iterative design loops.

A common bridge between molecular structure and sequence models is the Simplified Molecular-Input Line-Entry System (SMILES) (Weininger, 1988; Weininger et al., 1989), which linearizes a molecular graph into a string. This linearization has enabled a large family of sequence-based molecular methods, including (1) conditional SMILES generation from text prompts (Bagal et al., 2022; Zhumagambetov et al., 2021; Irwin et al., 2022), and (2) SMILES or structure-conditioned captioning that maps molecules to natural-language descriptions (Edwards et al., 2022; Zeng et al., 2024). Many of these systems rely on *autoregressive* (AR) pretrained language models (PLMs) such as GPT (Brown et al., 2020), T5 (Raffel et al., 2020), and BART (Lewis et al., 2020), which generate tokens sequentially by conditioning each next token on the prefix (Bagal et al., 2022; Zhumagambetov et al., 2021; Irwin et al., 2022). In this work, while we retain SMILES as the standard input/output interface, we explicitly project these molecules into a Knowledge Graph (KG) view represented as a serialized sequence of atom-bond-atom triplets to achieve better generalization and reduce dependency on specific SMILES syntax rules (e.g., ring number).

While effective, autoregressive models are not always well suited to SMILES-centric molecule–language modeling. First, decoding is inherently left-to-right, so early mistakes can propagate and are difficult to correct (Nie et al., 2025; Arriola et al., 2025). Second, SMILES encodes long-range syntactic and chemical dependencies (e.g., ring closures, branching, and valid substructures) that require coordinated decisions across distant positions; locally plausible token choices can therefore cascade into globally invalid molecules. Third, many objectives are inherently structure-level such as scaffold preservation that are hard to enforce

*Equal contribution

under strict prefix conditioning. These limitations motivate non-autoregressive generation mechanisms that can revise all positions jointly, raising a broader question: *can we build models that generate and explain molecules while explicitly preserving global validity and structural controllability?*

Diffusion language models offer a promising non-autoregressive alternative by iteratively denoising a corrupted representation, enabling holistic updates over all positions at each step and better coordination of long-range dependencies (Li et al., 2022; Nie et al., 2025). However, standard diffusion applies noise uniformly across tokens, which can corrupt chemically critical tokens (e.g., ring indices and branching markers) too early in the process and hinder recoverability.

To address this, we propose BiMol-Diff, a unified diffusion framework that addresses the paired tasks of $text \rightarrow SMILES$ generation and $SMILES \rightarrow text$ captioning under a common denoising formulation. Unlike standard diffusion, BiMol-Diff employs a token-aware noising strategy that assigns token-dependent noise levels using the per-token training loss as a proxy for recovery difficulty. Intuitively, tokens that are consistently harder to denoise receive a more conservative corruption schedule, improving recoverability. We apply this mechanism consistently across both directions: for $SMILES \rightarrow text$, it helps preserve semantic tokens during denoising, and for $text \rightarrow SMILES$, it helps preserve chemically salient SMILES tokens during generation.

Contributions. We make three primary contributions: (1) a token-aware noising strategy which preserves critical structural semantics during the diffusion process; (2) a unified framework that solves the bidirectional tasks of molecule generation and captioning, utilizing the same token-aware mechanism for both inverse problems; and (3) strong empirical performance across a wide range of metrics on both molecule generation and molecule captioning benchmarks.

2 Background and Preliminaries

2.1 Related Work

Diffusion for Molecular Graphs (Unconditional / Property-based): Diffusion models are widely used for molecular graph generation, spanning unconditional sampling and property-conditioned design. Recent variants improve expressivity via transformer-style graph denoisers (Liu et al.,

2024a), explore hierarchical/latent formulations for scalability and structure preservation (Bian et al., 2024), and study control mechanisms for property prediction and guidance (Zhang et al., 2025).

Text-guided Molecule Generation: Text-to-molecule generation is typically formulated as learning $p(\tilde{\mathcal{G}} | \mathbf{S})$ (or $p(\text{SMILES} | \mathbf{S})$), where \mathbf{S} describes desired structure or function. Recent work adapts PLMs to chemical strings, enabling $text \leftrightarrow molecule$ translation with encoder-decoder pretraining (Edwards et al., 2022) and improved multi-task / chemistry-aware extensions (Kim et al., 2025b). In parallel, diffusion-based text conditioning has emerged, including diffusion-LM style text-guided generation in discrete/embedded sequence spaces (Gong et al., 2024) and graph/latent diffusion variants that preserve molecular structure (Chang and Ye, 2025).

Molecule Captioning / Molecule-to-Text: Generating faithful descriptions from molecules (G2S) has been explored via translation-style baselines (Edwards et al., 2022) and molecule captioning models that incorporate structural encoders to improve grounding (Liu et al., 2023b). More recent multimodal approaches further strengthen molecule-language alignment for captioning and related tasks.

Unifying / Bidirectional Molecule-Text Modeling: While text-guided generation and captioning have progressed independently, *fully bidirectional* modeling that supports both $p(\mathbf{S} | \tilde{\mathcal{G}})$ and $p(\tilde{\mathcal{G}} | \mathbf{S})$ within a single generation framework remains relatively limited. Translation-based PLM systems offer practical bidirectionality (Edwards et al., 2022), and recent foundation-model efforts move toward broader cross-modal generalization (Liu et al., 2023b), but diffusion-based approaches are typically developed for one direction (most often S2G) and do not explicitly address token-level corruption heterogeneity. BiMol-Diff fills this gap by casting *both* directions under a unified conditional diffusion view and introducing a molecular graph-aware, token-dependent noising schedule to preserve chemically salient information during generation (see Appendix A.1).

2.2 Preliminaries

2.2.1 Diffusion Models: Forward and Reverse Process

Denoising diffusion probabilistic models (DDPMs) are generative models that learn a data distribution,

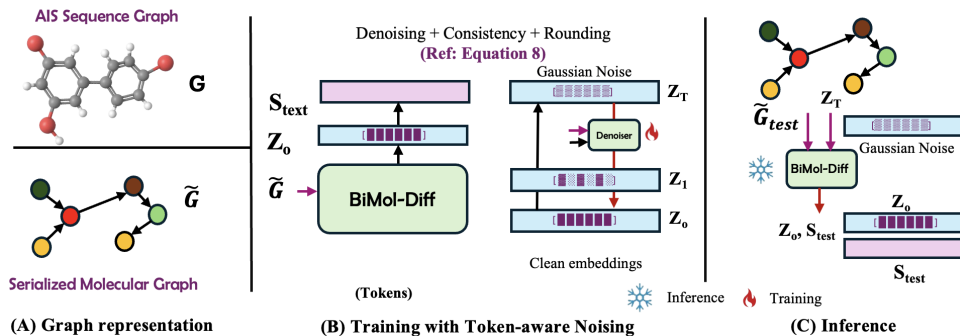


Figure 1: **The BiMol-Diff Framework.** (A) Molecules are represented through canonical SMILES, AIS-aware tokens, and serialized graph-triplet sequences. (B) **Training:** the model is trained with token-aware noising that preserves chemically salient tokens, optimizing denoising, consistency, and rounding objectives. (C) **Inference:** starting from Gaussian noise, iterative denoising and rounding generate either a caption (G2S) or a serialized molecular graph sequence (S2G), which is then deterministically decoded into canonical SMILES.

often conditioned on some context \mathbf{c} , $p(\mathbf{z}_0 | \mathbf{c})$. They consist of a fixed forward process and a learned reverse process.

Forward process: A standard DDPM forward process corrupts clean data \mathbf{z}_0 through a Markov chain with noise-schedule coefficients $\{\alpha_t\}_{t=1}^T$ controlling signal decay. This yields the standard closed-form for sampling a noised state \mathbf{z}_t at timestep t :

$$\mathbf{z}_t = \sqrt{\bar{\alpha}_t} \mathbf{z}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \quad (1)$$

with $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$ and $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Standard diffusion models typically use a fixed, data-agnostic (isotropic) noise schedule.

Reverse process with Conditional Denoising: The reverse process learns to recover the clean data \mathbf{z}_0 from pure noise $\mathbf{z}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. It is defined as a Markov chain $p_\theta(\mathbf{z}_{0:T})$ where each reverse transition $p_\theta(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{c})$ is a Gaussian whose mean $\boldsymbol{\mu}_\theta$ and variance $\boldsymbol{\Sigma}_\theta$ are parameterized by a model $\mathcal{M}_\theta(\mathbf{z}_t, t, \mathbf{c})$. The model is trained to predict the mean of the true posterior $q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{z}_0)$. The model parameters θ are optimized by maximizing the variational lower bound (VLB) on the conditional log-likelihood:

$$\begin{aligned} \mathcal{L}_{\text{vlb}} = \mathbb{E}_q \left[\underbrace{-\log p_\theta(\mathbf{z}_0 | \mathbf{z}_1, \mathbf{c})}_{\text{Reconstruction } (L_0)} \right. \\ \left. + \sum_{t=2}^T \underbrace{D_{\text{KL}}(q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{z}_0) \| p_\theta(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{c}))}_{\text{Denoising Matching } (L_{t-1})} \right. \\ \left. + \underbrace{D_{\text{KL}}(q(\mathbf{z}_T | \mathbf{z}_0) \| p(\mathbf{z}_T))}_{\text{Prior Matching } (L_T)} \right]. \quad (2) \end{aligned}$$

While this formulation is tractable, direct optimization of the full VLB is often unstable.

2.2.2 Molecules as Graphs:

A molecule is represented as a graph $\mathcal{G} = (\mathbf{V}, \mathbf{X}, \mathbf{A}, \mathbf{P})$, where $\mathbf{V} = \{1, \dots, n\}$ is the set of atoms (nodes) and $|\mathbf{V}| = n$. The atom feature matrix $\mathbf{X} \in \mathbb{R}^{n \times a}$ contains an a -dimensional feature vector for each atom. The adjacency tensor $\mathbf{A} \in \mathbb{R}^{n \times n \times b}$ encodes bond existence and bond types, where $\mathbf{A}_{ij} \in \{0, 1\}^b$ indicates the type of the bond between atoms i and j . The coordinate matrix $\mathbf{P} \in \mathbb{R}^{n \times 3}$ stores the 3D positions of atoms. Here, a is the atom feature dimension and b is the bond type. Depending on the task, subsets of these components may be used, e.g., $(\mathbf{V}, \mathbf{X}, \mathbf{A})$ for 2D graphs or $(\mathbf{V}, \mathbf{X}, \mathbf{P})$ for 3D conformations.

3 Methodology

3.1 Problem Statement

We use canonical SMILES as the standard molecule representation throughout BiMol-Diff. From this canonical string view, we derive two related representations. First, we denote by $G = \{m_1, \dots, m_K\}$ the Atoms-in-SMILES (AIS) (Ucak et al., 2023) token sequence, which provides a chemistry-aware alternative to character-level SMILES. Second, for graph-conditioned modeling, we denote by \tilde{G} the serialized molecular graph obtained by writing the molecule as an edge-list sequence of atom-bond-atom triplets using special tokens [HEAD], [REL], [TAIL], and [SEP]:

$$\tilde{G} = \langle [\text{HEAD}] h_i [\text{REL}] r_{ij} [\text{TAIL}] t_j [\text{SEP}] \rangle_{(i,j)},$$

where (h_i, r_{ij}, t_j) denotes an atom-bond-atom triple. Let $S_{\text{text}} = \{s_1, \dots, s_N\}$ be a textual description. We study two conditional generation tasks: (i) **molecule captioning** $\tilde{G} \rightarrow S_{\text{text}}$, and (ii)

molecule generation $S_{\text{text}} \rightarrow \tilde{G}$. Accordingly, we learn $M_\theta : \tilde{G} \rightarrow S_{\text{text}}$ and $M_\phi : S_{\text{text}} \rightarrow \tilde{G}$ and maximize

$$\max_{\Theta} \log p(S_{\text{text}} | \tilde{G}; \theta) + \log p(\tilde{G} | S_{\text{text}}; \phi),$$

where $\Theta = \{\theta, \phi\}$. We next describe the molecule-side representation pipeline used to realize these two directions, starting from canonical SMILES and proceeding through AIS-aware tokens and serialized graph sequences.

3.2 The BiMol-Diff Framework Overview

We introduce BiMol-Diff, a unified diffusion framework that addresses molecular graph-to-sequence (G2S) and sequence-to-graph (S2G) generation under a single conditional modeling view. Given paired data $(\mathcal{A}, \mathcal{B}) \in \{(S_{\text{text}}, \tilde{G}), (\tilde{G}, S_{\text{text}})\}$, BiMol-Diff learns the conditional distribution $p(\mathcal{A} | \mathcal{B})$, where \mathcal{B} provides the conditioning context and \mathcal{A} is the target modality.

Molecule Representation Pipeline. We use canonical SMILES as the common molecule interface in both directions and convert it to its Atoms-in-SMILES representation G (Ucak et al., 2023). For graph-conditioned modeling, this molecule-side representation is deterministically mapped to the serialized graph sequence \tilde{G} , where each bond is written as an atom-bond-atom triplet using [HEAD]/[REL]/[TAIL] templates and [SEP] delimiters. Thus, in the G2S direction, \tilde{G} serves as the conditioning sequence for generating S_{text} . Conversely, in the S2G direction, S_{text} conditions prediction of \tilde{G} , and the predicted \tilde{G} is deterministically mapped back through the same molecule-side representation pipeline into canonical SMILES for evaluation. Figure 2 summarizes this symmetric encoding/decoding view.

Conditional Diffusion. For either direction, we embed the target \mathcal{A} into continuous “clean” latents $\mathbf{z}_0 = g_\Phi(\mathcal{A})$ and apply a conditional DDPM to obtain \mathbf{z}_t . The reverse processes denoise conditioned on the source modality: $\mathcal{M}_\theta(\mathbf{z}_t, t, \tilde{G})$ for G2S and $\mathcal{M}_\phi(\mathbf{z}_t, t, S_{\text{text}})$ for S2G, enabling a unified objective across both directions.

Token-aware Noising. Unlike standard data-agnostic schedules that corrupt all tokens uniformly, BiMol-Diff uses a token-aware, token-dependent noising strategy that preserves chemically salient tokens more conservatively. We detail the construction and use of this schedule in

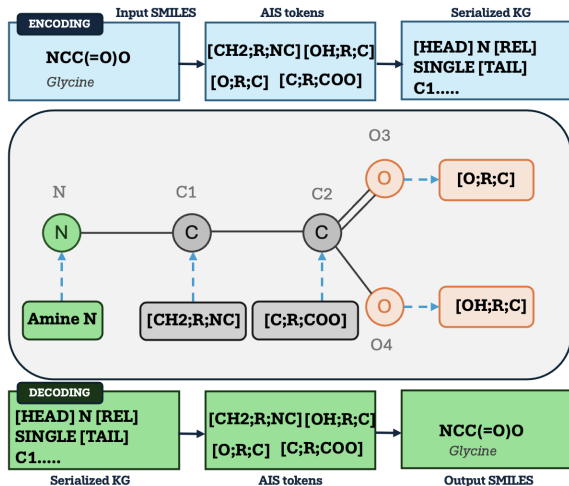


Figure 2: Molecule encoding and decoding in BiMol-Diff.

Sec 3.5. This improves recoverability in both directions while keeping the same diffusion formulation across the unified framework.

3.3 Graph Encoding for BiMol-Diff

In Section 3.1, a molecular graph is a set of relational triplets $\tilde{G} = \{(h_i, r_{ij}, t_j)\}$. To integrate with transformer based encoder-decoder backbone, we serialize this set of triplets into a single token sequence. Following the PLM-based G2S convention, each triplet is mapped to a short template $\langle [\text{HEAD}] h_i [\text{REL}] r_{ij} [\text{TAIL}] t_j \rangle$, and all such segments are concatenated with a separator token [SEP]:

$$\tilde{G} = \langle [\text{HEAD}] h_{i_1} [\text{REL}] r_{i_1 j_1} [\text{TAIL}] t_{j_1} \rangle [\text{SEP}] \dots [\text{SEP}] \langle [\text{HEAD}] h_{i_M} [\text{REL}] r_{i_M j_M} [\text{TAIL}] t_{j_M} \rangle \quad (3)$$

where $\{(i_m, j_m)\}_{m=1}^M$ enumerates all bonds with $i_m < j_m$ to avoid duplicates. This linearized sequence \tilde{G} serves as the conditioning signal to BiMol-Diff in the graph-to-sequence task. Figure 2 (top) illustrates this encoding process.

Example (molecule KG \rightarrow caption). Consider the molecule with SMILES string CCO (ethanol). Its molecular graph has three atoms $h_1 = C_1$, $h_2 = C_2$, $h_3 = O_3$ and two single bonds: (1, 2) and (2, 3). The corresponding KG triplets are $(C_1, \text{SINGLE}, C_2)$ and $(C_2, \text{SINGLE}, O_3)$. Our serialized KG becomes

$$\tilde{G} = \langle [\text{HEAD}] C_1 [\text{REL}] \text{SINGLE} [\text{TAIL}] C_2 \rangle [\text{SEP}] \langle [\text{HEAD}] C_2 [\text{REL}] \text{SINGLE} [\text{TAIL}] O_3 \rangle$$

In the G2S setting, BiMol-Diff takes the graph-triplet representation \tilde{G} as input and generates

a natural-language caption \mathbf{S}_{text} describing the molecule. The reverse realization path used in the S2G direction is described separately in Section 3.4.

3.4 Graph Decoding for BiMol-Diff

In the S2G direction, conditioned on \mathbf{S}_{text} , BiMol-Diff predicts the serialized molecular graph sequence \tilde{G} . We decode this sequence by first parsing it into atom-bond-atom triplets (h_i, r_{ij}, t_j) , and then merging these triplets to recover the corresponding molecular graph. The recovered graph is subsequently mapped to its AIS token representation G and canonicalized into a standard SMILES string for S2G evaluation. Figure 2 (bottom) illustrates this decoding path.

Example (caption \rightarrow molecule KG). Consider the textual description corresponding to ethanol. In the S2G direction, BiMol-Diff predicts the serialized graph sequence:

$$\tilde{G} = \langle [\text{HEAD}] C_1 [\text{REL}] \text{SINGLE} [\text{TAIL}] C_2 \rangle [\text{SEP}] \langle [\text{HEAD}] C_2 [\text{REL}] \text{SINGLE} [\text{TAIL}] O_3 \rangle. \quad (4)$$

These triplets are merged to recover the molecular graph with atoms C_1, C_2, O_3 and bonds (1, 2) and (2, 3). The recovered graph is then mapped to its molecule-side representation and canonicalized into the SMILES string CCO.

3.5 Token Aware Noising

Motivation & Rationale: Standard diffusion models rely on fixed, data-agnostic noising schedules that apply noise uniformly across all tokens.

We argue this is suboptimal for molecular tasks because chemically salient tokens (e.g., functional groups, stereochemistry) typically carry higher information density than structurally redundant or low-information tokens.

Under uniform noising, chemically salient tokens and simple syntactic tokens are corrupted at the same rate, which can weaken molecular fidelity. To address this, we introduce a *Token-aware noise schedule*. We hypothesize that the model’s per-timestep reconstruction loss ℓ_t^i can serve as a useful proxy for token difficulty and importance. By learning a mapping from ℓ_t^i to the corresponding token’s noise schedule, we preserve high information tokens while allowing information-redundant tokens to tolerate higher noise levels. This re-parametrizes the denoising path, aiming to improve generation quality in both G2S and S2G tasks.

Noising Schedule: Following the rationale discussed in Section 3.5, for the G2S setting, we propose a token-dependent noising schedule, parameterized by a vector $\bar{\alpha}_{t,\text{new}}^i \in \mathbb{R}^N$, unlike the uniform noising with baseline cumulative schedule $\bar{\alpha}_t$ (for e.g.: *sqrt*) used by standard diffusion models. This formulation has two stages summarized in Algo 1. **Stage 1: Estimating token-wise difficulty.** For each token i and diffusion step $t = 1, \dots, T$, we define the denoising difficulty as:

$$\ell_t^i = \mathbb{E}_{\mathbf{z}_t \sim q(\mathbf{z}_t | \mathbf{z}_0)} \|\mathcal{M}_\theta(\mathbf{z}_t, t, \tilde{G})^{(i)} - \mathbf{z}_0^{(i)}\|^2 \quad (5)$$

Averaging over the training set yields a difficulty profile $(\ell_1^i, \dots, \ell_T^i)$ for each i . Empirically, ℓ_t^i tends to increase with t (later steps are noisier), but the estimated profile is not strictly monotone. We also compute $\ell_{\min}^i = \min_t \ell_t^i$ and $\ell_{\max}^i = \max_t \ell_t^i$ to define the difficulty range for token i . In Stage 2 these profiles and their ranges are used to construct a token-wise cumulative schedule, and to obtain a monotonic difficulty profile for each token i .

Stage 2: Token-aware schedule. Given $(\ell_t^i)_{t=1}^T$ and the baseline cumulative schedule $(\bar{\alpha}_t)_{t=1}^T$, we construct an adaptive schedule $(\bar{\alpha}_{t,\text{new}}^i)_{t=1}^T$ for each token i . Since this schedule controls noise applied at each step, we want to reallocate noise according to denoising difficulty. Hence, we define a piecewise-linear map $\Psi_i : [\ell_{\min}^i, \ell_{\max}^i] \rightarrow (0, 1)$ that interpolates the baseline schedule as a function of loss:

$$\Psi_i(x) = \bar{\alpha}_{t-1} + \frac{\bar{\alpha}_t - \bar{\alpha}_{t-1}}{\ell_t^i - \ell_{t-1}^i} (x - \ell_{t-1}^i), \quad (6)$$

with $x \in [\ell_{t-1}^i, \ell_t^i]$, $t = 2, \dots, T$, $\Psi_i(\ell_1^i) = \bar{\alpha}_1$ and $\Psi_i(\ell_T^i) = \bar{\alpha}_T$. In case $\ell_t^i = \ell_{t-1}^i$, we add a tiny jitter ε to avoid division by zero. Empirically, $(\ell_t^i)_{t=1}^T$ is not strictly monotone in t , so instead of using its raw values we introduce a new linear ramp in difficulty space:

$$\ell_t^{i,\text{new}} = \ell_{\min}^i + \frac{t-1}{T-1} (\ell_{\max}^i - \ell_{\min}^i) \quad (7)$$

with $t = 1, \dots, T$. Substituting this $\ell_t^{i,\text{new}}$ into $\Psi_i(x)$ we get a new cumulative schedule $\bar{\alpha}_{t,\text{new}}^i = \Psi_i(\ell_t^{i,\text{new}})$ for $t = 1, \dots, T$. We clamp $\bar{\alpha}_{t,\text{new}}^i$ to $(0, 1)$ and apply a non-increasing isotonic projection (refer Appendix A.3) over t to obtain the final schedule $0 < \bar{\alpha}_{t+1,\text{new}}^i \leq \bar{\alpha}_{t,\text{new}}^i < 1$ for all t .

Method	#P	B	CrF++	M	B-F1	MVE
Autoregressive Baselines						
MolT5-Base (Edwards et al., 2022)	220M	0.452	0.651	0.510	0.681	0.852
Text+Chem T5 (Christofidellis et al., 2023)	223M	<u>0.542</u>	<u>0.701</u>	0.648	0.728	0.866
MolCA (Liu et al., 2023b)	110M	0.531	0.665	0.651	0.709	0.815
GitMol (Liu et al., 2024b)	700M	0.475	0.680	0.532	0.751	0.875
GraphT5 (Kim et al., 2025a)	272M	0.481	0.692	0.545	<u>0.810</u>	<u>0.913</u>
Diffusion Baselines						
Diffusion-LM (Li et al., 2022)	91M	0.512	0.702	0.602	0.783	0.861
DiffuSeq (Gong et al., 2023)	91M	0.532	0.708	0.601	0.812	0.887
TGM-DLM (Gong et al., 2024)	125M	0.467	0.689	0.589	0.779	0.856
BiMol-Diff (ours)	63M	0.567	0.734	0.626	0.843	0.925
%Gain (vs. Best AR)	x3.5↓	+4.6%	+4.7%	-3.8%	+4.1%	+1.3%
%Gain (vs. Best Diff)	x1.4↓	+6.6%	+3.7%	+4.2%	+3.8%	+4.3%

Table 1: Molecule captioning performance on the M3-20M dataset. Benchmarking BiMol-Diff against SoTA Autoregressive (AR) and Diffusion models. BiMol-Diff achieves state-of-the-art results, surpassing the best diffusion and AR baselines by substantial margins across all metrics while using fewer parameters.

Algorithm 1 Token Aware Noise Schedule

Require: Baseline cumulative schedule $\{\tilde{\alpha}_t\}_{t=1}^T$, update interval K

Ensure: Schedules $\{\tilde{\alpha}_{t,\text{new}}^i\}_{t=1}^T$ for all i

- 1: **if** train_step % $K == 0$ **then**
- 2: **for** all i **do**
- 3: Estimate $\{\ell_t^i\}_{t=1}^T$ via Eq. 5; compute $\{\ell_{\min}^i, \ell_{\max}^i\}$.
- 4: Define piecewise-linear map Ψ_i as in Eq. (6).
- 5: Construct difficulty ramp $\{\ell_t^{i,\text{new}}\}_{t=1}^T$ via Eq. (7).
- 6: Compute new schedule $\tilde{\alpha}_t^i = \Psi_i(\ell_t^{i,\text{new}})$.
- 7: Clamp $\tilde{\alpha}_t^i$ to $(0, 1)$, apply a non-increasing isotonic projection to obtain $\{\tilde{\alpha}_{t,\text{new}}^i\}_{t=1}^T$.
- 8: **end for**
- 9: **end if**
- 10: **return** $\{\tilde{\alpha}_{t,\text{new}}^i\}$ for all tokens i and steps t .

3.6 Model Training and Inference

Training: Our training objective is derived from the Variational Lower Bound (VLB) (Eq. 2) presented in the preliminaries. While the full VLB optimization can be unstable (Ho et al., 2020), a common simplification is to train the model \mathcal{M}_θ to predict the added noise ϵ . However, our framework adopts an alternative \mathbf{z}_0 -prediction reparameterization, which trains the model to directly predict the clean data \mathbf{z}_0 at every timestep t . A critical component of this objective is the rounding term $L_0 = -\log \tilde{p}_\Phi(\mathbf{S}_{\text{text}} | \mathbf{z}_0)$, which handles the final step of converting the continuous latent

variable \mathbf{z}_0 back into discrete tokens \mathbf{S}_{text} . We define this as a trainable rounding distribution: $\tilde{p}_\Phi(\mathbf{S}_{\text{text}} | \mathbf{z}_0) = \prod_{i=1}^N \tilde{p}_\Phi(s_i | \mathbf{z}_{0,i})$, where each token s_i is sampled from a softmax distribution over the vocabulary, using logits derived from the corresponding output embedding $\mathbf{z}_{0,i}$. By combining this rounding term (for $t = 0$) with the denoising matching terms (for $t > 1$) using our \mathbf{z}_0 -prediction reparameterization, we arrive at our final, composite objective. (The full derivation from the VLB is in App A.2).

$$\mathcal{L}_{\text{e2e-simple}}(\mathbf{S}_{\text{text}}) = \mathbb{E}_q \left[\underbrace{\sum_{t=2}^T \|\mathcal{M}_\theta(\mathbf{z}_t, t, \tilde{\mathcal{G}}) - \mathbf{z}_0\|^2}_{\text{Denoising}} + \underbrace{\|g_\Phi(\mathbf{S}_{\text{text}}) - \mathcal{M}_\theta(\mathbf{z}_1, 1, \tilde{\mathcal{G}})\|^2}_{\text{Consistency}} - \underbrace{\log \tilde{p}_\Phi(\mathbf{S}_{\text{text}} | \mathbf{z}_0)}_{\text{Rounding}} \right] \quad (8)$$

The same objective is used in the S2G direction, with the target discrete sequence replaced by $\tilde{\mathcal{G}}$ and the conditioning sequence replaced by S_{text} . This objective directly optimizes the most critical parts of the process: the denoising accuracy across all steps (Denoising), the consistency of the first denoising step with the true data embedding (Consistency), and the quality of the final conversion to discrete tokens (Rounding).

Inference-time schedule: The \mathbf{z}_0 -prediction objective (Eq: 8) trains the denoiser to explicitly predict the clean latent at each step. At inference, given \mathbf{z}_t and conditioning context \mathbf{c} , we compute $\hat{\mathbf{z}}_0 = \mathcal{M}_\theta(\mathbf{z}_t, t, \mathbf{c})$ (with $\mathbf{c} = \tilde{\mathcal{G}}$ for G2S and $\mathbf{c} = \mathbf{S}_{\text{text}}$ for S2G). We then use our learned token-wise cumulative schedule $\tilde{\alpha}_{t,\text{new}}^i$ in the reverse sampling update, where $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. To mitigate rounding errors (refer section 3.6), we apply the clamping trick, which forces the predicted

vector to commit to a word for intermediate diffusion steps:

$$\mathbf{z}_{t-1}^{(i)} = \sqrt{\bar{\alpha}_{t-1,\text{new}}^i} \text{Clamp}\left(\mathcal{M}_\theta(\mathbf{z}_t, t, \mathbf{c})^{(i)}\right) + \sqrt{1 - \bar{\alpha}_{t-1,\text{new}}^i} \epsilon^{(i)}, \quad (9)$$

where $\text{Clamp}(\cdot)$ denotes nearest-neighbor projection onto the embedding table. This strategy ensures that the diffusion trajectory remains grounded in valid token embeddings during sampling.

4 Experiments

4.1 Experimental Setup

Model Architecture. We implement BiMol-Diff using two distinct encoder-decoder Transformer models tailored for each task utilizing GeLU activations. Crucially, both tasks are trained for 200,000 steps using our proposed *Token-aware noise schedule* (refer §3.5) over $T = 2000$ diffusion steps. We adapt the architecture for each task as follows: for Molecule Generation (S2G), we employ a 6-encoder/12-decoder configuration ($d_{\text{model}} = 1024$; $\approx 180\text{M}$ parameters). The encoder processes frozen SciBERT embeddings, while the decoder predicts the serialized molecular graph sequence \tilde{G} . This predicted graph sequence is then deterministically mapped back into a canonical SMILES string for evaluation. This variant is trained with a batch size of 64 and a learning rate of 5×10^{-5} . Conversely, for Molecule Captioning (G2S), we utilize a 6-encoder/9-decoder configuration ($d_{\text{model}} = 512$; $\approx 63\text{M}$ parameters) with a peak learning rate of 10^{-4} . This model encodes serialized molecular structure using the AIS-based SMILES vocabulary augmented with learnable special tokens ([HEAD], [REL], [TAIL], [SEP]), while the decoder generates captions utilizing the bert-base-uncased vocabulary (Devlin et al., 2019).

Datasets. We evaluate our framework on two distinct datasets tailored to each task. For Molecule Generation (S2G), we utilize the ChEBI-20 (Edwards et al., 2022) benchmark, comprising 33,010 molecules with SMILES strings restricted to a maximum length of 256. This dataset is partitioned into an 80%/10%/10% split for training, validation, and testing, respectively. For Molecule Captioning (G2S), we employ a filtered subset of the M3-20M dataset (Guo et al., 2025b), containing approximately 360,000 SMILES-description pairs. This subset, also divided into an 80%/10%/10% split,

pairs molecular graphs with diverse textual descriptions derived from curated scientific literature or synthesized via GPT-3.5.

Baselines. We compare BiMol-Diff against state-of-the-art (SoTA) autoregressive and diffusion-based methods for each task. For Molecule Generation (S2G), we benchmark against MolXPT (Liu et al., 2023a) as a competitive autoregressive SMILES generator, and leading diffusion baselines including 3M-Diffusion (Zhu et al., 2024) and UT-GDiff (Xiang et al., 2025). For Molecule Captioning (G2S), we compare against MolT5-Base (Edwards et al., 2022) as a standard SMILES→text sequence-to-sequence baseline, GraphT5 (Kim et al., 2025a) as a recent graph-text model, and DiffuSeq (Gong et al., 2023) as a representative diffusion-based sequence generator.

Evaluation Metrics. For Molecule Captioning (G2S), we evaluate the quality of generated text using both surface-level and semantic measures. We report standard n -gram overlap metrics including BLEU (Papineni et al., 2002), METEOR (Banerjee and Lavie, 2005), and ChrF++ (Popović, 2015), alongside embedding-based metrics such as BERTScore-F1 (Zhang* et al., 2020) and MAUVE (Pillutla et al., 2023) to assess semantic alignment and distributional divergence. For Molecule Generation (S2G), we focus on structural fidelity and chemical plausibility. We report Exact Match accuracy to measure precise reconstruction of the ground-truth molecule. To assess the capture of chemical substructures, we compute structural overlap using Tanimoto similarity on MACCS, RDKit, and Morgan fingerprints. Finally, we report chemical validity to ensure the generated SMILES strings correspond to valid molecular graphs.

4.2 Experimental Results

Molecule Captioning (G2S): Table 1 reports results on M3-20M. BiMol-Diff achieves the best overall performance on BLEU and ChrF++, reaching **0.567** BLEU and **0.734** ChrF++, which improves over the best performing autoregressive baseline (Text+Chem T5) by **4.6%** and **4.7%**, respectively. On semantic metrics, our model attains **0.843** BERTScore-F1 and **0.925** MAUVE, outperforming graph-aware baselines such as GraphT5, and consistently improving over diffusion models (e.g., DiffuSeq) across all reported metrics. While MolCA achieves the highest METEOR score, BiMol-Diff remains competitive (**0.626**) and offers a strong overall trade-off across sur-

Method	Type	#P	Exact	MACCS	RDKit	Morgan	Valid
Text-Based Autoregressive (SMILES)							
T5-Base (Raffel et al., 2020)	AR	248M	0.069	0.731	0.605	0.545	0.660
MolT5-Base (Edwards et al., 2022)	AR	248M	0.081	0.721	0.588	0.529	0.772
MolXPT (Liu et al., 2023a)	AR	350M	0.215	0.859	0.757	0.667	0.983
T5 Enc + MolXPT MLP (Deng et al., 2025)	Adapter	111M	0.227	0.820	0.713	0.636	0.980
Text-Based Autoregressive (SELFIES)							
SciBERT + MolGen (Deng et al., 2025)	Adapter	317M	0.083	0.680	0.526	0.422	0.995
Galactica-1.3B + MolGen (Fang et al., 2024)	Adapter	1.5B	0.091	0.706	0.560	0.454	<u>0.995</u>
T5 Enc + MolGen (Deng et al., 2025)	Adapter	317M	0.165	0.758	0.616	0.527	0.995
Graph & Diffusion Baselines							
DiGress (sim. guidance) (Vignac et al., 2023)	Diff	289M	0.014	0.577	0.389	0.288	0.854
3M-Diffusion (Zhu et al., 2024)	Diff	162M	0.005	0.548	0.370	0.273	1.000
Graph-DiT (Liu et al., 2024a)	Diff	162M	0.000	0.374	0.269	0.159	0.909
UTGDiff (w/o pretrain) (Xiang et al., 2025)	Diff	125M	<u>0.227</u>	<u>0.867</u>	<u>0.763</u>	<u>0.695</u>	0.856
BiMol-Diff (ours)	Diff	180M	0.262	0.894	0.791	0.762	0.901
%Gain (vs. Best AR)	–	x1.9↓	+15.4%	+4.1%	+4.5%	+14.2%	-8.3%
%Gain (vs. Best Diff)	–	x1↓	+15.4%	+3.1%	+3.7%	+9.6%	+5.3%

Table 2: Molecule generation performance on the ChEBI-20 test set. BiMol-Diff outperforms state-of-the-art Autoregressive and Diffusion baselines in structure reconstruction (Exact Match) and fingerprint similarity, validating the effectiveness of token-aware denoising.

face and embedding-based measures, supporting the benefit of token-aware noising.

Molecule Generation (S2G): Table 2 summarizes results on ChEBI-20. BiMol-Diff attains an Exact Match of **0.262**, improving over both the best autoregressive adapter baseline (T5 Enc + MolXPT MLP) and the strongest diffusion baseline (UTGDiff), each at 0.227 (+15.4% relative). This gain is reflected in fingerprint similarities, achieving **0.894** (MACCS), **0.791** (RDKit), and **0.762** (Morgan), including a **14.2%** improvement in Morgan similarity over the best autoregressive baseline. Finally, BiMol-Diff maintains high chemical validity (**0.901**) while prioritizing reconstruction accuracy, whereas some baselines attain higher validity with weaker structural recovery.

4.3 Ablation

We analyze three aspects of BiMol-Diff: (i) the proposed token-aware noising strategy and its mapping function, (ii) the effect of molecular tokenization, and (iii) decoding efficiency of BiMol-Diff. **(i) Impact of noise schedule and choice of mapping function.** Table 3 (top) shows that replacing a uniform corruption schedule (*sqr*t) with our token-aware schedule yields a substantial improvement (0.495→0.567 BLEU; 0.682→0.734 ChrF++; 0.531→0.626 METEOR). Among mapping choices, the linear mapping performs best, indicating that explicitly allocating lower corruption to harder-to-recover tokens is beneficial for caption recovery. Figure 3 visualizes how token-aware schedules differ from a global *sqr*t schedule. Instead of applying the same noise profile to all positions, we learn per-token schedules that modulate

the corruption level according to token difficulty: positions with higher denoising loss receive more conservative corruption (lower noise at intermediate timesteps), while easier positions can be noised more aggressively. This selective corruption better preserves content during the forward process and improves recoverability in the reverse process, which aligns with the gains observed in the main results (Table 1, 2).

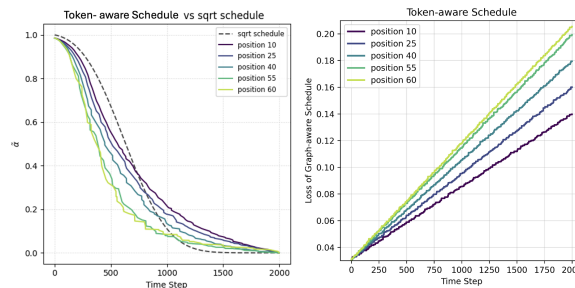


Figure 3: Token-aware noising vs. uniform *sqr*t schedule (captioning). *Left*: the learned token-wise schedules apply non-uniform corruption across positions compared to a global *sqr*t schedule. *Right*: token-wise difficulty profiles are used to construct schedules that apply more conservative corruption to harder-to-recover tokens.

(ii) Impact of molecular tokenization. Table 3 (bottom) applies the best noising configuration (token-aware + linear) with different molecular tokenizers. Our current framework uses an AIS-based SMILES tokenizer for a stable, deterministic segmentation that preserves SMILES syntax (e.g., bracketed atoms and ring/branch symbols) and enables straightforward serialization, the tokenizer is plug-and-play and can be swapped without changing the denoising objective or architecture. We observe a clear monotonic trend: regex-based

Noise Schedule	Mapping / Tokenizer	B	CrF++	M
<i>Impact of Noise Schedule & Mapping Function</i>				
Uniform	Sqrt Schedule	0.495	0.682	0.531
Token-Aware	Cosine Mapping	0.548	0.720	0.595
Token-Aware (Ours)	Linear Mapping	0.567	0.734	0.626
<i>Impact of Molecular Tokenization</i>				
Token-Aware (Linear)	Regex-based	0.558	0.714	0.606
Token-Aware (Linear)	Atom-level	0.563	0.722	0.613
Token-Aware (Linear)	Atoms-in-SMILES (AIS)	0.567	0.734	0.626

Table 3: Ablation study on the M3-20M dataset. We analyze the impact of different (**Top**) Noise Schedules & Mapping Functions and (**Bottom**) Tokenization strategies on molecule captioning performance. BiMol-Diff (Token-Aware + Linear) with Atoms-in-SMILES (AIS) tokenization yields the best results.

SMILES tokenization underperforms atom-level tokens, while Atoms-in-SMILES (AIS) yields the strongest performance across all reported metrics.

(iii) Decoding Efficiency. Given the iterative nature of diffusion models, we examine the trade-off between generation quality and inference latency. Reducing the number of reverse steps leads to substantial speedups, albeit with a degradation in output fidelity. To focus on modeling performance, all primary results are reported with $T=2000$ steps. Table 4 presents the quality-latency trade-off for Graph-to-Sequence tasks. Within the diffusion family, BiMol-Diff exhibits markedly improved efficiency over DiffuSeq. At $T=2000$, it achieves a BLEU score of 0.567 in 89 seconds per batch, yielding a $3.6\times$ speedup over DiffuSeq (317 seconds), while also improving generation quality. When compared to Autoregressive (AR) baselines,

Method	Params	Steps	Time (s)	Speedup	BLEU \uparrow
<i>Autoregressive Baselines</i>					
MolT5-Base (Edwards et al., 2022)	220M	–	<5s	–	0.452
Text+Chem T5 (Christofidellis et al., 2023)	223M	–	<5s	–	0.542
<i>Diffusion Baselines</i>					
DiffuSeq (Gong et al., 2023)	91M	2000	317s	1.0 \times (Ref)	0.532
<i>Ours</i>					
BiMol-Diff	63M	2000	89s	3.6 \times	0.567
BiMol-Diff	63M	1000	45s	7.0 \times	0.551
BiMol-Diff	63M	500	23s	13.8 \times	0.365
BiMol-Diff	63M	100	5s	63.4 \times	0.312

Table 4: Inference efficiency comparison on the G2S task. We compare BiMol-Diff against the best Autoregressive (AR) and Diffusion baselines from Table 1. Time is measured as total inference time for a batch size of 50 on a single NVIDIA V100 GPU.

our model shows a clear operational boundary. In the high-quality output range (1000–2000 steps), BiMol-Diff outperforms the strongest AR baseline, Text+Chem T5 (0.542 BLEU), confirming that the extra computational time yields better text generation. However, the model struggles when pushed for extreme speed. As we reduce the infer-

ence steps below 1000, we observe a sharp drop in performance. At 500 steps, although the inference time improves to 23 seconds ($13.8\times$ speedup), the BLEU score falls to 0.365, which is lower than even the base MolT5 model (0.452). This indicates that BiMol-Diff requires a minimum threshold of denoising steps (approximately 1000) to maintain its advantage. Our model is ideal for high-precision generation and simple AR models remain the better choice for low-latency applications.

5 Conclusion

We introduced BiMol-Diff, a unified diffusion framework addressing the bidirectional tasks of molecular graph generation and textual captioning. By moving beyond standard data-agnostic corruption to a *token-aware noise schedule*, our approach explicitly preserves chemically salient substructures during the diffusion process. Empirical validation on M3-20M and ChEBI-20 confirms that this strategy yields significant gains over state-of-the-art autoregressive and diffusion baselines, particularly in structural fidelity (S2G) and semantic alignment (G2S). These findings demonstrate that token-aware denoising is essential for high-fidelity, controllable molecular language modeling, offering a robust alternative to autoregressive approaches in scientific domains.

6 Limitations

First, BiMol-Diff relies on linearized sequence representations and standard attention, lacking the explicit structural inductive bias inherent to graphs, which may limit data efficiency for complex topologies. Second, as with most diffusion models, our approach incurs higher computational costs during inference compared to autoregressive baselines due to the iterative denoising process, highlighting the need for future work on accelerated sampling.

References

- Marianne Arriola, Subham Sekhar Sahoo, Aaron Gokaslan, Zhihan Yang, Zhixuan Qi, Jiaqi Han, Justin T Chiu, and Volodymyr Kuleshov. 2025. [Block diffusion: Interpolating between autoregressive and diffusion language models](#). In *The Thirteenth International Conference on Learning Representations*.
- Viraj Bagal, Rishal Aggarwal, P. K. Vinod, and U. Deva Priyakumar. 2022. [Molgpt: Molecular generation using a transformer-decoder model](#). *Journal of Chemical Information and Modeling*, 62(9):2064–2076. PMID: 34694798.
- Satanjeev Banerjee and Alon Lavie. 2005. [METEOR: An automatic metric for MT evaluation with improved correlation with human judgments](#). In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.
- Tian Bian, Yifan Niu, Heng Chang, Divin Yan, Junzhou Huang, Yu Rong, Tingyang Xu, Jia Li, and Hong Cheng. 2024. [Hierarchical graph latent diffusion model for conditional molecule generation](#). In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management, CIKM '24*, page 130–140, New York, NY, USA. Association for Computing Machinery.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, and 12 others. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Jinho Chang and Jong Chul Ye. 2025. [LDMol: A text-to-molecule diffusion model with structurally informative latent space surpasses AR models](#). In *Forty-second International Conference on Machine Learning*.
- Dimitrios Christofidellis, Giorgio Giannone, Jannis Born, Ole Winther, Teodoro Laino, and Matteo Manica. 2023. [Unifying molecular and textual representations via multi-task language modelling](#). In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 6140–6157. PMLR.
- Yifan Deng, Spencer S. Ericksen, and Anthony Gitter. 2025. [Chemical language model linker: Blending text and molecules with modular adapters](#). *Journal of Chemical Information and Modeling*, 65(17):8944–8956.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Carl Edwards, Tuan Lai, Kevin Ros, Garrett Honke, Kyunghyun Cho, and Heng Ji. 2022. [Translation between molecules and natural language](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 375–413, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Yin Fang, Ningyu Zhang, Zhuo Chen, Lingbing Guo, Xiaohui Fan, and Huajun Chen. 2024. [Domain-agnostic molecular generation with chemical feedback](#). In *The Twelfth International Conference on Learning Representations*.
- Bahare Fatemi, Jonathan Halcrow, and Bryan Perozzi. 2024. [Talk like a graph: Encoding graphs for large language models](#). In *The Twelfth International Conference on Learning Representations*.
- Haisong Gong, Qiang Liu, Shu Wu, and Liang Wang. 2024. [Text-guided molecule generation with diffusion language model](#). In *Proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence and Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence and Fourteenth Symposium on Educational Advances in Artificial Intelligence, AAAI'24/IAAI'24/EAAI'24*. AAAI Press.
- Shansan Gong, Mukai Li, Jiangtao Feng, Zhiyong Wu, and Lingpeng Kong. 2023. [DiffuSeq: Sequence to sequence text generation with diffusion models](#). In *The Eleventh International Conference on Learning Representations*.
- Shuhan Guo, Yatao Bian, Ruibing Wang, Nan Yin, Zhen Wang, and Quanming Yao. 2025a. [Unified molecule-text language model with discrete token representation](#). In *Proceedings of the Thirty-Fourth International Joint Conference on Artificial Intelligence, IJCAI-25*, pages 9205–9213. International Joint Conferences on Artificial Intelligence Organization. AI4Tech: AI Enabling Technologies.
- Siyuan Guo, Lexuan Wang, Chang Jin, Jinxian Wang, Han Peng, Huayang Shi, Wengen Li, Jihong Guan, and Shuigeng Zhou. 2025b. [M³-20m: A large-scale multi-modal molecule dataset for ai-driven drug design and discovery](#). *Journal of Bioinformatics and Computational Biology*, 23(02):2550006. PMID: 40494666.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. [Denosing diffusion probabilistic models](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 6840–6851. Curran Associates, Inc.

- Ross Irwin, Spyridon Dimitriadis, Jiazhen He, and Esben Jannik Bjerrum. 2022. [Chemformer: a pre-trained transformer for computational chemistry](#). *Machine Learning: Science and Technology*, 3(1):015022.
- Sangyeup Kim, Nayeon Kim, Yinhua Piao, and Sun Kim. 2025a. [GraphT5: Unified molecular graph-language modeling via multi-modal cross-token attention](#). *Preprint*, arXiv:2503.07655.
- Seojin Kim, Hyeontae Song, Jaehyun Nam, and Jinwoo Shin. 2025b. [Training text-to-molecule models with context-aware tokenization](#). In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 22442–22460, Suzhou, China. Association for Computational Linguistics.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Xiang Lisa Li, John Thickstun, Ishaan Gulrajani, Percy Liang, and Tatsunori B. Hashimoto. 2022. [Diffusion-lm improves controllable text generation](#). In *Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS '22*, Red Hook, NY, USA. Curran Associates Inc.
- Gang Liu, Jiaxin Xu, Tengfei Luo, and Meng Jiang. 2024a. [Graph diffusion transformers for multi-conditional molecular generation](#). In *Advances in Neural Information Processing Systems*. Curran Associates, Inc.
- Pengfei Liu, Yiming Ren, Jun Tao, and Zhixiang Ren. 2024b. [GIT-Mol: A multi-modal large language model for molecular science with graph, image, and text](#). *Computers in Biology and Medicine*, 171:108073.
- Zequan Liu, Wei Zhang, Yingce Xia, Lijun Wu, Shufang Xie, Tao Qin, Ming Zhang, and Tie-Yan Liu. 2023a. [MolXPT: Wrapping molecules with text for generative pre-training](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1606–1616, Toronto, Canada. Association for Computational Linguistics.
- Zhiyuan Liu, Sihang Li, Yanchen Luo, Hao Fei, Yixin Cao, Kenji Kawaguchi, Xiang Wang, and Tat-Seng Chua. 2023b. [MolCA: Molecular graph-language modeling with cross-modal projector and uni-modal adapter](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 15623–15638, Singapore. Association for Computational Linguistics.
- Shen Nie, Fengqi Zhu, Zebin You, Xiaolu Zhang, Jingyang Ou, Jun Hu, Jun Zhou, Yankai Lin, Ji-Rong Wen, and Chongxuan Li. 2025. [Large language diffusion models](#). *Preprint*, arXiv:2502.09992.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: a method for automatic evaluation of machine translation](#). In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, PA. Association for Computational Linguistics.
- Krishna Pillutla, Lang Liu, John Thickstun, Sean Welleck, Swabha Swayamdipta, Rowan Zellers, Sewoong Oh, Yejin Choi, and Zaid Harchaoui. 2023. [Mauve scores for generative models: Theory and practice](#). *Preprint*, arXiv:2212.14578.
- Maja Popović. 2015. [chrF: character n-gram F-score for automatic MT evaluation](#). In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *Journal of Machine Learning Research*, 21(140):1–67.
- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. 2015. [Deep unsupervised learning using nonequilibrium thermodynamics](#). In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 2256–2265, Lille, France. PMLR.
- Duong Thanh Tran, Nhat Truong Pham, Nguyen Doan Hieu Nguyen, and Balachandran Manavalan. 2024. [Mol2Lang-VLM: Vision- and text-guided generative pre-trained language models for advancing molecule captioning through multimodal fusion](#). In *Proceedings of the 1st Workshop on Language + Molecules (L+M 2024)*, pages 97–102, Bangkok, Thailand. Association for Computational Linguistics.
- Umit Ucak, Islambek Ashyrmamatov, and Juyong Lee. 2023. [Improving the quality of chemical language model outcomes with atom-in-smiles tokenization](#). *Journal of Cheminformatics*, 15.
- Clement Vignac, Igor Krawczuk, Antoine Siraudin, Bohan Wang, Volkan Cevher, and Pascal Frossard. 2023. [DiGress: Discrete denoising diffusion for graph generation](#). In *The Eleventh International Conference on Learning Representations*.
- David Weininger. 1988. [Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules](#). *Journal of Chemical Information and Computer Sciences*, 28(1):31–36.

David Weininger, Arthur Weininger, and Joseph L. Weininger. 1989. [Smiles. 2. algorithm for generation of unique smiles notation](#). *Journal of Chemical Information and Computer Sciences*, 29(2):97–101.

Yuran Xiang, Haiteng Zhao, Chang Ma, and Zhi-Hong Deng. 2025. [Instruction-based molecular graph generation with unified text-graph diffusion model](#). *Preprint*, arXiv:2408.09896.

Zheni Zeng, Bangchen Yin, Shipeng Wang, Jiarui Liu, Cheng Yang, Haishen Yao, Xingzhi Sun, Maosong Sun, Guotong Xie, and Zhiyuan Liu. 2024. [ChatMol: interactive molecular discovery with natural language](#). *Bioinformatics*, 40(9):btae534.

Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. 2020. [BERTScore: Evaluating text generation with bert](#). In *International Conference on Learning Representations*.

Xiaochen Zhang, Shuangxi Wang, Ying Fang, and Qiankun Zhang. 2025. [MG-DIFF: A novel molecular graph diffusion model for molecular generation and optimization](#). *PLOS ONE*, 20(10):e0331450.

Zaixiang Zheng, Yifan Deng, Dongyu Xue, Yi Zhou, Fei Ye, and Quanquan Gu. 2023. [Structure-informed language models are protein designers](#). In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 42317–42338. PMLR.

Huasheng Zhu, Teng Xiao, and Vasant G Honavar. 2024. [3M-Diffusion: Latent multi-modal diffusion for language-guided molecular structure generation](#). In *First Conference on Language Modeling*.

Rustam Zhumagambetov, Ferdinand Molnar, Vsevolod A. Peshkov, and Siamac Fazli. 2021. [Transmol: repurposing a language model for molecular generation](#). *RSC Adv.*, 11:25921–25932.

A Appendix

A.1 Related Works

Table 5 positions BiMol-Diff relative to recent state-of-the-art molecule-text generation frameworks. While autoregressive (AR) baselines, such as MolT5 and UniMoT, have successfully demonstrated bidirectional capabilities, they often suffer from the standard limitations of sequential decoding (e.g., exposure bias). Conversely, existing diffusion-based approaches like TGM-DLM and UTGDiff offer non-autoregressive benefits but are predominantly unidirectional, focusing almost exclusively on the Text→Graph modality. BiMol-Diff bridges this gap as a unified diffusion framework capable of both Graph→Text and Text→Graph generation. Furthermore, we distinguish our approach by integrating a *token-aware*

mechanism, a feature notably absent in prior diffusion baselines—which adapts the noising schedule to prioritize chemically significant tokens, similar to the strategy employed by the AR-based ChemT5.

A.2 Derivation

BiMol-Diff builds on the standard diffusion framework, which trades the flexibility of expressive generative models (e.g., GANs, VAEs, flow models) for the tractability of likelihood-based training in a continuous latent space \mathbf{z} . The overall goal is to minimize the negative log-likelihood

$$\mathbb{E}_{\mathbf{z}_0, \mathbf{c}}[-\log p_\theta(\mathbf{z}_0 | \mathbf{c})], \quad (10)$$

which is upper-bounded by the Variational Lower Bound (VLB).

A.2.1 Forward and Reverse Processes

The forward Markov chain is defined as $q(\mathbf{z}_{1:T} | \mathbf{z}_0) = \prod_{t=1}^T q(\mathbf{z}_t | \mathbf{z}_{t-1})$, where each transition is Gaussian:

$$q(\mathbf{z}_t | \mathbf{z}_{t-1}) = \mathcal{N}\left(\mathbf{z}_t | \sqrt{1 - \beta_t} \mathbf{z}_{t-1}, \beta_t \mathbf{I}\right). \quad (11)$$

Let $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$. By induction, the marginal at time t satisfies:

$$\mathbf{z}_t = \sqrt{\bar{\alpha}_t} \mathbf{z}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (12)$$

so that $q(\mathbf{z}_t | \mathbf{z}_0) = \mathcal{N}\left(\sqrt{\bar{\alpha}_t} \mathbf{z}_0, (1 - \bar{\alpha}_t) \mathbf{I}\right)$. We used the *sqr*t schedule as the baseline schedule used in DiffusionLM (Li et al., 2022), namely $\bar{\alpha}_t = 1 - \sqrt{t/T + s}$ with small $s > 0$. The reverse denoising process then learns

$$p_\theta(\mathbf{z}_{0:T}) = p(\mathbf{z}_T) \prod_{t=1}^T p_\theta(\mathbf{z}_{t-1} | \mathbf{z}_t), \quad (13)$$

$$p_\theta(\mathbf{z}_{t-1} | \mathbf{z}_t) = \mathcal{N}\left(\boldsymbol{\mu}_\theta(\mathbf{z}_t, t), \boldsymbol{\sigma}_\theta^2(\mathbf{z}_t, t)\right).$$

Applying Bayes’ rule to the forward transitions yields the exact posterior mean

$$\boldsymbol{\mu}_t(\mathbf{z}_t, \mathbf{z}_0) = \frac{\sqrt{\bar{\alpha}_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{z}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \mathbf{z}_0, \quad (14)$$

whose coefficients we denote by \mathcal{U} and \mathcal{E} . BiMol-Diff’s training objective is then to match the network’s predicted $\boldsymbol{\mu}_\theta, \boldsymbol{\sigma}_\theta$ to these posterior quantities via a simple noise-prediction loss. We optimize the negative log-likelihood by upper-bounding it with the variational lower bound

$$\mathbb{E}[-\log p_\theta(x_0)] \leq \mathcal{L}_{\text{vlb}} = \sum_{t=0}^T \mathcal{L}_t. \quad (15)$$

Method	Generation family	Text-conditioned	Graph→Text	Text→Graph	Bidirectional	Token-aware training
<i>Autoregressive (AR)</i>						
MolT5 (Edwards et al., 2022)	Seq2Seq PLM (T5)	✓	✓	✓	✓	—
Text+ChemT5 (Christofidellis et al., 2023)	Seq2Seq PLM (T5)	✓	✓	✓	✓	—
CAMT5 (Kim et al., 2025b)	Seq2Seq PLM (T5)	✓	—	✓	—	✓
UniMoT (Guo et al., 2025a)	LLM-based multitask	✓	✓	✓	✓	—
LDMol (Chang and Ye, 2025)	PLM-based generator	✓	—	✓	—	—
MolCA (Liu et al., 2023b)	Encoder-decoder (graph-aware)	—	✓	—	—	—
Mol2Lang-VLM (Tran et al., 2024)	Multimodal PLM	—	✓	—	—	—
<i>Non-autoregressive (Diffusion)</i>						
TGM-DLM (Gong et al., 2024)	Diffusion-LM	✓	—	✓	—	—
3M-Diffusion (Zhu et al., 2024)	Latent diffusion (graph/seq)	✓	—	✓	—	—
UTGDiff (Xiang et al., 2025)	Graph diffusion	✓	—	✓	—	—
BiMol-Diff (ours)	Conditional diffusion	✓	✓	✓	✓	✓

Table 5: Feature comparison of BiMol-Diff against representative autoregressive (AR) and diffusion baselines. (✓ = supported; — = not supported/reported). Token-aware mechanism denotes chemistry-aware tokenization/weighting in AR models and token-wise noising schedules in diffusion models. Unlike prior diffusion works which are typically unidirectional, BiMol-Diff enables fully bidirectional generation.

A.2.2 Variational Lower Bound (VLB)

Following Sohl-Dickstein et al. (Sohl-Dickstein et al., 2015), for conditional generation the VLB decomposes into:

$$\begin{aligned} \mathcal{L}_{\text{vlb}} &= \mathbb{E}_{q(\mathbf{z}_{1:T}|\mathbf{z}_0)} \left[\mathcal{L}_T + \sum_{t=2}^T \mathcal{L}_t - \mathcal{L}_0 \right], \\ \mathcal{L}_T &:= \log \frac{q(\mathbf{z}_T | \mathbf{z}_0)}{p(\mathbf{z}_T)}, \\ \mathcal{L}_t &:= \log \frac{q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{z}_0)}{p_\theta(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{c})}, \\ \mathcal{L}_0 &:= \log p_\theta(\mathbf{z}_0 | \mathbf{z}_1, \mathbf{c}). \end{aligned} \quad (16)$$

where each \mathcal{L}_t is a KL divergence between Gaussians. The true posterior mean (via Bayes’ rule) is:

$$\boldsymbol{\mu}_t(\mathbf{z}_t, \mathbf{z}_0) = \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}}_{\mathcal{U}} \mathbf{z}_t + \underbrace{\frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t}}_{\mathcal{E}} \mathbf{z}_0, \quad (17)$$

with covariance $\boldsymbol{\Sigma}_q = \tilde{\beta}_t \mathbf{I}$, $\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$. In the standard simplification, the model’s covariance is fixed to match the true posterior covariance $\boldsymbol{\Sigma}_\theta = \boldsymbol{\Sigma}_q$, the KL collapses to a weighted MSE:

$$\mathcal{L}_t = \frac{1}{2} \left\| \boldsymbol{\mu}_t - \boldsymbol{\mu}_\theta \right\|_{\boldsymbol{\Sigma}_q^{-1}}^2 \propto \mathbb{E} \left[\left\| \mathbf{z}_0 - \mathcal{M}_\theta(\mathbf{z}_t, t, \mathbf{c}) \right\|^2 \right]. \quad (18)$$

Thus, for $2 \leq t \leq T$, $\mathcal{L}_t \rightarrow \left\| \mathbf{z}_0 - \mathcal{M}_\theta(\mathbf{z}_t, t, \mathbf{c}) \right\|^2$. The final KL encourages \mathbf{z}_T to match the unit Gaussian prior:

$$\mathcal{L}_T = \text{KL}(q(\mathbf{z}_T | \mathbf{z}_0) \| p(\mathbf{z}_T)) \propto \left\| \boldsymbol{\mu}(\mathbf{z}_T) \right\|^2, \quad (19)$$

a constant w.r.t. θ . The discrete target \mathbf{S} (sequence) is encoded into a continuous embedding $g_\Phi(\mathbf{S})$.

The final term in VLB is $\mathcal{L}_0 = -\log p_\theta(\mathbf{z}_0 | \mathbf{z}_1, \mathbf{c})$. We need to integrate the discrete data \mathbf{S} into this continuous likelihood term. We use the law of total probability to express the continuous likelihood $p_\theta(\mathbf{z}_0 | \mathbf{z}_1, \mathbf{c})$ by marginalizing over all possible discrete tokens in the target sequence $\mathbf{S} = \{s_1, s_2, \dots, s_N\}$:

$$p_\theta(\mathbf{z}_0 | \mathbf{z}_1, \mathbf{c}) = \sum_{\mathbf{S}} p_\theta(\mathbf{z}_0, \mathbf{S} | \mathbf{z}_1, \mathbf{c}) \quad (20)$$

We then apply the product rule to the joint probability:

$$p_\theta(\mathbf{z}_0, \mathbf{S} | \mathbf{z}_1, \mathbf{c}) = p_\theta(\mathbf{z}_0 | \mathbf{S}, \mathbf{z}_1, \mathbf{c}) \cdot p_\theta(\mathbf{S} | \mathbf{z}_1, \mathbf{c}) \quad (21)$$

For training, we are interested in the specific ground-truth sequence \mathbf{S} . When we evaluate \mathcal{L}_0 during training, we consider only the term where \mathbf{S} is the ground-truth sequence:

$$\begin{aligned} \mathcal{L}_0 &\approx -\log p_\theta(\mathbf{z}_0, \mathbf{S} | \mathbf{z}_1, \mathbf{c}) \\ &= -\log \left[p_\theta(\mathbf{z}_0 | \mathbf{S}, \mathbf{z}_1, \mathbf{c}) \cdot p_\theta(\mathbf{S} | \mathbf{z}_1, \mathbf{c}) \right]. \end{aligned} \quad (22)$$

The core approximation simplifies the dependency graph by asserting that the discrete data \mathbf{S} is generated only from the clean latent \mathbf{z}_0 , and is independent of \mathbf{z}_1 and \mathbf{c} given \mathbf{z}_0 .

$$\mathbf{S} \perp (\mathbf{z}_1, \mathbf{c}) | \mathbf{z}_0 \quad (23)$$

This allows us to replace the discrete conditional likelihood with the separate rounding network $\tilde{p}_\Phi(\mathbf{S} | \mathbf{z}_0)$: $p_\theta(\mathbf{S} | \mathbf{z}_1, \mathbf{c}) \approx \tilde{p}_\Phi(\mathbf{S} | \mathbf{z}_0)$. Substituting this back into the likelihood decomposition:

$$p_\theta(\mathbf{z}_0, \mathbf{S} | \mathbf{z}_1, \mathbf{c}) \approx p_{\text{cont}}(\mathbf{z}_0 | \mathbf{S}, \mathbf{z}_1, \mathbf{c}) \cdot \tilde{p}_\Phi(\mathbf{S} | \mathbf{z}_0) \quad (24)$$

Taking the negative logarithm of the approximation gives the two desired terms:

$$\mathcal{L}_0 \approx -\log p_{\text{cont}}(\mathbf{z}_0 \mid \mathbf{S}, \mathbf{z}_1, \mathbf{c}) - \log \tilde{p}_\Phi(\mathbf{S} \mid \mathbf{z}_0)$$

This split yields the two components used in the final training objective:

1. Consistency Term ($\mathcal{L}_{\text{Cons}}$): The first term is the negative log-likelihood of the continuous latent, which is minimized via the MSE loss on the means: $-\log p_{\text{cont}}(\mathbf{z}_0 \mid \mathbf{S}, \mathbf{z}_1, \mathbf{c}) \rightarrow \mathcal{L}_{\text{Consistency}} = \left\| g_\Phi(\mathbf{S}) - \mathcal{M}_\theta(\mathbf{z}_1, 1, \mathbf{c}) \right\|^2$.
2. Rounding Term ($\mathcal{L}_{\text{Round}}$): This second term is the dedicated loss for the discrete data likelihood: $\mathcal{L}_{\text{Round}} = -\log \tilde{p}_\Phi(\mathbf{S} \mid \mathbf{z}_0)$

A.2.3 Final End-to-End Objective

Combining all components:

$$\mathcal{L}_{\text{vbl}} \propto \sum_{t=2}^T \underbrace{\left\| \mathbf{z}_0 - \mathcal{M}_\theta(\mathbf{z}_t, t, \mathbf{c}) \right\|^2}_{\text{Denoising}} \quad (25)$$

$$+ \underbrace{\left\| g_\Phi(\mathbf{S}) - \mathcal{M}_\theta(\mathbf{z}_1, 1, \mathbf{c}) \right\|^2}_{\text{Consistency}} \quad (26)$$

$$- \underbrace{\log \tilde{p}_\Phi(\mathbf{S} \mid \mathbf{z}_0)}_{\text{Rounding}}. \quad (27)$$

Dropping constant terms, the simplified end-to-end training loss is:

$$\begin{aligned} \mathcal{L}_{\text{e2e-simple}}(\mathbf{S}) = \mathbb{E}_q \left[\sum_{t=2}^T \underbrace{\left\| \mathcal{M}_\theta(\mathbf{z}_t, t, \tilde{\mathcal{G}}) - \mathbf{z}_0 \right\|^2}_{\text{Denoising}} \right. \\ \left. + \underbrace{\left\| g_\Phi(\mathbf{S}) - \mathcal{M}_\theta(\mathbf{z}_1, 1, \tilde{\mathcal{G}}) \right\|^2}_{\text{Consistency}} \right. \\ \left. - \underbrace{\log \tilde{p}_\Phi(\mathbf{S} \mid \mathbf{z}_0)}_{\text{Rounding}} \right]. \end{aligned} \quad (28)$$

A.3 Non-increasing Isotonic Projection

After constructing the per-token cumulative schedule $\tilde{\alpha}_t^i$, we project it onto the set of non-increasing sequences $\{\bar{\alpha}_t^i\}_{t=1}^T$ such that $\bar{\alpha}_1^i \geq \bar{\alpha}_2^i \geq \dots \geq \bar{\alpha}_T^i$. Concretely, this is a 1D isotonic regression problem with squared loss, which we solve using the standard Pool-Adjacent-Violators Algorithm (PAVA). This algorithm finds the closest monotone non-increasing sequence (in the least-squares sense)

to the input. Intuitively, it smooths out spurious "bumps" in the loss profile while guaranteeing that the cumulative signal strength strictly decays over time, fulfilling the monotonicity requirement of the diffusion process.