

Bridging the Discrete-Continuous Gap: Unified Multimodal Generation via Coupled Manifold Discrete Absorbing Diffusion

Anonymous ACL submission

Abstract

The bifurcation of generative modeling into autoregressive approaches for discrete data (text) and diffusion approaches for continuous data (images) hinders the development of truly unified multimodal systems. While Masked Language Models (MLMs) offer efficient bidirectional context, they traditionally lack the generative fidelity of autoregressive models and the semantic continuity of diffusion models. Furthermore, extending masked generation to multimodal settings introduces severe alignment challenges and training instability. In this work, we propose **CoM-DAD** (Coupled Manifold Discrete Absorbing Diffusion), a novel probabilistic framework that reformulates multimodal generation as a hierarchical dual-process. CoM-DAD decouples high-level semantic planning from low-level token synthesis. First, we model the semantic manifold via a continuous latent diffusion process; second, we treat token generation as a discrete absorbing diffusion process, regulated by a **Variable-Rate Noise Schedule**, conditioned on these evolving semantic priors. Crucially, we introduce a **Stochastic Mixed-Modal Transport** strategy that aligns disparate modalities without requiring heavy contrastive dual-encoders. Our method demonstrates superior stability over standard masked modeling, establishing a new paradigm for scalable, unified text-image generation.

1 Introduction

The pursuit of Artificial General Intelligence (AGI) necessitates models capable of reasoning and generating across diverse modalities. However, a fundamental topological disconnect persists in current architectures: language is inherently discrete and symbolic, while visual data is continuous and dense. Consequently, the field has fragmented into two dominant paradigms: Autoregressive (AR) models, which excel at discrete text generation, and Contin-

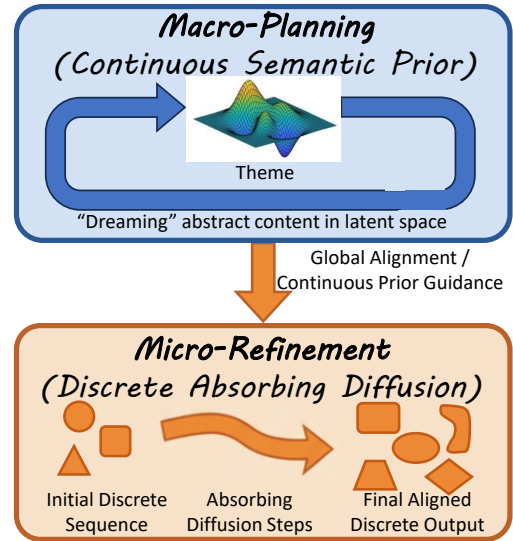


Figure 1: **Overview of CoM-DAD.** The framework splits generation into **Macro-Planning** (Top), where a continuous latent diffusion models abstract semantic "themes" (the "dreaming" phase), and **Micro-Refinement** (Bottom), where a discrete absorbing diffusion synthesizes tokens. The vertical arrow signifies the conditioning of the discrete generation process on the continuous prior, ensuring global alignment between the abstract plan and the final token sequence.

uous Diffusion Models (CDMs), which dominate high-fidelity image synthesis.

Efforts to unify these paradigms often result in compromised hybrids. Masked Generative Models (MGMs) attempt to bridge this gap by treating generation as a parallel denoising task. While MGMs offer significant efficiency gains over AR models (which suffer from serial dependency) and faster inference than CDMs, they notoriously struggle with **generative consistency**. Without a strong prior, masking-based models often produce locally coherent but globally disjoint outputs. Recent advancements, such as Representation-Conditioned Generation (RCG) (Li et al., 2024), have mitigated this in the visual domain by conditioning pixel generation on self-supervised representations. However, RCG remains strictly unimodal, focusing on

060 unconditional image synthesis and failing to ad- 110
061 dress the complexities of cross-modal alignment or 111
062 the discrete nature of language. 112

063 We argue that the difficulty in training multi- 113
064 modal masked models stems from forcing a single 114
065 network to simultaneously learn semantic abstrac- 115
066 tion (what to generate) and structural composition 116
067 (how to arrange tokens). To resolve this, we in- 117
068 troduce CoM-DAD, a hierarchical framework that
069 mathematically formalizes masked generation not
070 as simple “filling in the blanks,” but as a Discrete
071 Absorbing Diffusion Process guided by a Continu-
072 ous Semantic Prior. As conceptually illustrated in
073 Figure 1, this hierarchical decoupling allows us to
074 manage semantic abstraction and structural com-
075 position on separate, optimized manifolds. Our
076 approach operates on two levels:

- 077 1. Macro-Planning (Latent Space): We employ 123
078 a lightweight diffusion model to navigate the 124
079 continuous manifold of high-level semantic 125
080 representations. This allows the model to 126
081 “dream” the abstract content of an image or 127
082 sentence before committing to specific tokens. 128
- 083 2. Micro-Refinement (Discrete Space): We for- 129
084 mulate token generation as a reverse diffu- 130
085 sion process where tokens emerge from an 131
086 absorbing state ([MASK]). Unlike standard 132
087 MLMs, our transition kernel is strictly con- 133
088 ditioned on the macro-plan, ensuring that ev- 134
089 ery generated token is globally aligned with 135
090 the intended semantic target. Crucially, this 136
091 process is governed by a Variable-Rate Noise 137
092 Schedule, which replaces fixed masking ra- 138
093 tios with a continuous time parameter. This 139
094 allows the model to learn generation from 140
095 pure noise, while the transition kernel remains 141
096 strictly conditioned on the macro-plan to en- 142
097 sure global alignment. 143
098

099 Furthermore, to address the scarcity of aligned 144
100 multimodal data, we propose a **Stochastic Mixed-** 145
101 **Modal Transport** mechanism during training. By 146
102 dynamically swapping semantic priors between 147
103 modalities (e.g., forcing the model to generate im- 148
104 age tokens from a text representation), we induce a 149
105 unified semantic space without the need for auxil-
106 iary alignment losses like CLIP.

107 Our contributions can be summarized as follows:

- 108 • We propose **CoM-DAD**, a unified probabilis- 150
109 tic framework that bridges topological gaps 151

110 between modalities by coupling a continuous 111
112 latent diffusion for semantic planning with a 113
114 discrete absorbing diffusion for synthesis. 115

- We introduce a **Variable-Rate Discrete Diffu-** 113
114 **sion** mechanism that generalizes masked mod- 114
115 eling with a continuous noise schedule, signif- 115
116 icantly improving generative consistency over 116
117 static masking strategies. 117
- We develop a **Stochastic Mixed-Modal** 118
119 **Transport** strategy that naturally aligns visual 119
120 and textual manifolds via dynamic representa- 120
121 tion swapping, eliminating the need for heavy 121
122 contrastive dual-encoder pre-training. 122
- We demonstrate that our hierarchical decou- 123
124 pling achieves superior training stability and 123
125 sampling efficiency compared to standard au- 124
126 toregressive or monolithic diffusion baselines 125
127 in multimodal contexts. 126

128 2 Related Work

Diffusion Models for Discrete Generation. Dif- 129
130 fusion models have achieved notable success in 130
131 continuous domains such as image and audio gener- 131
132 ation (Sohl-Dickstein et al., 2015; Ho et al., 2020). 132
133 Extending diffusion to discrete sequences has at- 133
134 tracted increasing interest (Li et al., 2025), particu- 134
135 larly for language and symbolic reasoning. Early 135
136 works (Hoogeboom et al., 2021; Austin et al., 2021) 136
137 adapt diffusion to discrete spaces via relaxation or 137
138 masking strategies. Later approaches (Li et al., 138
139 2022; He et al., 2022) embed discrete tokens into 139
140 continuous spaces to enable Gaussian diffusion. 140
141 Although effective, this introduces optimization 141
142 difficulties, as simultaneously optimizing both the 142
143 embedding layer and the diffusion model can lead 143
144 to shortcut learning. In contrast, **CoM-DAD** op- 144
145 erates directly on the discrete manifold. By introduc- 145
146 ing a hierarchical coupling with a continuous latent 146
147 planner, it enables stable, efficient, and controllable 147
148 generation without the optimization instability as- 148
149 sociated with joint embedding-diffusion training. 149

Masked Language Models as Discrete Ab- 150
151 **sorbing Processes.** Masked language modeling 151
152 (MLM) has recently been reinterpreted as a form 152
153 of discrete diffusion with absorbing states (Google 153
154 DeepMind, 2025; Nie et al., 2025; Wu et al., 2025; 154
155 Ye et al., 2025). While these methods effectively 155
156 approximate autoregressive generation through it- 156
157 erative unmasking, they typically restrict both se- 157

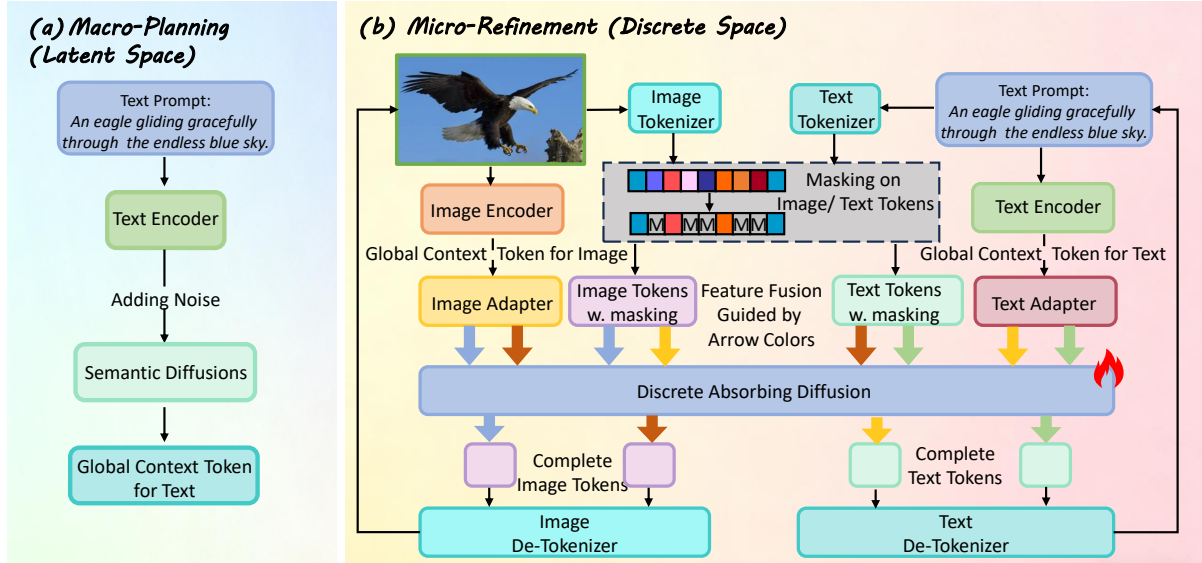


Figure 2: **The CoM-DAD Training Pipeline.** The framework consists of two coupled diffusion processes. **Left (Stage I):** The *Manifold-Constrained Semantic Diffusion* learns a continuous prior over semantic representations (r) via an SDE, capable of handling both text and image modalities. **Right (Stage II):** The *Semantic-Aware Discrete Absorbing Diffusion* reconstructs the discrete token sequence (x) from a masked state (\tilde{x}_t). The **Semantic Injection Interface** (center) connects these topologies by projecting the sampled semantic plan r into the decoder’s embedding space, conditioning the reverse diffusion step $p_\theta(x_{t-1}|\tilde{x}_t, r)$ to ensure global semantic coherence. Cross-Modal Alignment is applied to (b).

158 mantic abstraction and structural composition to
 159 the discrete token space. **CoM-DAD** distinguishes
 160 itself by hierarchically decoupling these tasks. We
 161 employ a continuous latent diffusion to generate
 162 a high-level semantic plan, which serves as a ro-
 163 bust condition for the discrete absorbing process,
 164 effectively bridging the topological gap between
 165 continuous semantic planning and discrete struc-
 166 tural generation.

167 **Multimodal Generation and Semantic Guidance.**

168 Multimodal generation aims to produce coherent
 169 outputs across heterogeneous data types. Prior
 170 work often relies on joint embedding spaces or
 171 autoregressive multimodal transformers (Alayrac
 172 et al., 2022; Chen et al., 2023; Team et al., 2023)
 173 to bridge the modality gap. **CoM-DAD** advances this
 174 paradigm by introducing Stochastic Mixed-Modal
 175 Transport. Rather than treating modalities as dis-
 176 parate sources to be aligned, we unify them into a
 177 shared continuous semantic manifold. This allows
 178 for dynamic prior swapping, where the discrete dif-
 179 fusion process is universally guided by the contin-
 180 uous planner, facilitating seamless cross-modal ge-
 181 neralization. Semantic-Aware generation (Li et al.,
 182 2024; Wang and Torr, 2022) demonstrates that high-
 183 level representations can effectively guide contin-
 184 uous diffusion, though existing methods remain
 185 unimodal. **CoM-DAD** extends this idea to multi-
 186 modal discrete generation by injecting learned rep-

187 resentations directly into the absorbing diffusion
 188 process, enabling efficient semantic planning and
 189 cross-modal alignment within a unified framework.

190 **3 Method**

191 In this section, we formally detail **CoM-DAD**, a
 192 hierarchical generative framework that bridges the
 193 gap between continuous semantic exploration and
 194 discrete token generation. Unlike prior semantic-
 195 conditioned approaches such as RCG (Li et al.,
 196 2024), which focus exclusively on unimodal image
 197 synthesis using standard backbones, **CoM-DAD** in-
 198 troduces a **unified discrete diffusion mechanism**
 199 capable of joint text-image modeling.

200 Our framework, schematically illustrated in Fig-
 201 ure 2, decomposes the intractable multimodal distri-
 202 bution $p(x)$ into two tractable generative processes
 203 operating in distinct topological spaces: (1) a **Con-**
 204 **tinuous Latent Diffusion** process modeling the
 205 high-level semantic manifold \mathcal{R} , and (2) a **Discrete**
 206 **Absorbing Diffusion** process modeling the token-
 207 space conditional distribution $p(x|r)$. The figure
 208 highlights how the Semantic Injection Interface
 209 acts as the critical bridge, translating the contin-
 210 uous "plan" from the latent diffusion into a guiding
 211 signal for the discrete token denoiser.

212 **3.1 Theoretical Formulation**

213 Let $x \in \mathcal{X}$ represent a discrete sequence (e.g., text
 214 tokens or quantized image patches) and $r \in \mathbb{R}^d$

Table 1: **Quantitative comparison of unconditional text generation performance.** CoM-DAD outperforms existing autoregressive and masked language model baselines in both BLEU-2 and BLEU-4 metrics. This superior fidelity validates the effectiveness of the Continuous Latent Planner in maintaining global coherence across the discrete text manifold. “Ours (large) + Autoregressive” denotes a variant where CoM-DAD is constrained to generate tokens in a sequential order.

Methods	Pretained	Training Steps	Inference Iterations	Output Length	BLEU /% (\uparrow)	Self-BLEU /% (\downarrow)
<i>Autoregressive Prediction</i>						
GPT-2	✓	1M	40	40	10.81	40.02
BERT (base) (Devlin, 2018)	✓	1M	40	40	7.80	10.06
BERT (large) (Devlin, 2018)	✓	1M	40	40	5.05	9.43
Ours (base) + Autoregressive	×	400K+300K	20	256	8.52	18.37
Ours (large) + Autoregressive	×	400K+300K	20	256	13.64	16.52
<i>Diffusion-based Methods</i>						
D3PM (Austin et al., 2021)	×	1M	128	128	42.41	22.88
Diffusion-LM (Li et al., 2022)	×	740K	2000	64	35.53	26.68
DiffusionBERT (He et al., 2022)	×	1.9M	128	128	43.58	21.51
BERT-Mouth (Wang and Cho, 2019)	✓	2.8M	10	50	28.67	12.4
Ours (base)	×	400K+300K	20	256	29.42	18.37
Ours (large)	×	400K+300K	20	256	47.46	16.52

be a continuous semantic vector derived from a pre-trained encoder $\mathcal{E}(x)$. We maximize the evidence lower bound (ELBO) of the log-likelihood $\log p_\theta(x)$:

$$\log p_\theta(x) \geq \underbrace{\mathbb{E}q(r|x)[\log p_\theta(x|r)]}_{\text{reconstruction}} - \underbrace{\text{DKL}(q(r|x)|p_\phi(r))}_{\text{prior matching}}. \quad (1)$$

Unlike standard VAEs where the prior $p(r)$ is a static Gaussian, we parameterize $p_\phi(r)$ as a **continuous diffusion model**. Furthermore, we model the reconstruction term $p_\theta(x|r)$ not as a simple autoregressive decoder, but as a **discrete diffusion process** over the vocabulary set V . This hybrid formulation allows CoM-DAD to decouple global semantic planning (in continuous space \mathcal{R}) from local structural refinement (in discrete space \mathcal{X}).

3.2 Stage I: Manifold-Constrained Semantic Diffusion

The first stage models the prior distribution of semantic representations $p_\phi(r)$. While RCG (Li et al., 2024) utilizes a standard diffusion model for image class embeddings, we employ a modality-agnostic diffusion process capable of navigating the joint semantic space of both vision and language.

Given a semantic vector $r_0 = \mathcal{E}(x)$, we define a forward stochastic differential equation (SDE) that gradually destroys semantic information:

$$dr_t = -\frac{1}{2}\beta(t)r_t dt + \sqrt{\beta(t)}d\mathbf{w}_t, \quad (2)$$

where \mathbf{w}_t is standard Brownian motion. We train a time-dependent denoiser $\epsilon_\phi(r_t, t)$ to reverse this process. Crucially, to ensure stability across modal-

ities with varying norms, we employ **representation normalization** before the diffusion process. The training objective is the standard reweighted variational bound:

$$\mathcal{L}_{\text{latent}} = \mathbb{E}_{t, r_0, \epsilon} [\|\epsilon - \epsilon_\phi(r_t, t, c_m)\|^2], \quad (3)$$

where c_m indicates the modality ID, allowing the single model to learn the distinct manifolds of textual (r_{txt}) and visual (r_{img}) semantics simultaneously.

3.3 Stage II: Semantic-Aware Discrete Absorbing Diffusion

The core innovation of CoM-DAD is the formulation of sequence generation as a **Discrete Absorbing Diffusion Process**, generalizing the concept of “masked language modeling” into a rigorous probabilistic framework.

Discrete Forward Process. Let $x_0 = (w_1, \dots, w_L)$ be a sequence of tokens from vocabulary V . We define a forward transition matrix Q_t that transitions any token w to a special absorbing state [MASK] with probability γ_t , and leaves it unchanged with probability $1 - \gamma_t$. This defines a corrupted sequence \tilde{x}_t where a subset of tokens are masked according to the Markov property of the absorbing state.

Variable-Rate Noise Schedule. A critical component of our approach is the noise schedule $\gamma(t)$. Unlike the fixed masking strategies used in standard BERT-like models (typically 15%), we sample the masking ratio γ_t from a continuous time

schedule $t \sim U[0, 1]$. This creates a **Variable-Rate Noise Schedule** that forces the model to learn generation across the entire complexity spectrum—from pure noise ($\gamma_1 \approx 1$) to fine-grained refinement ($\gamma_0 \approx 0$). This dynamic schedule is what effectively shifts the model’s capability from simple local infilling to robust global generation.

Semantic-Aware Denoising. We learn a reverse transition kernel $p_\theta(x_{t-1}|\tilde{x}_t, r)$ parameterized by a Transformer. To condition the discrete generation on the continuous semantic vector r sampled from Stage I, we introduce a **Semantic Injection Interface**:

$$h_0 = [\text{Proj}(r); \text{Embed}(\tilde{x}_t)]. \quad (4)$$

By projecting r into the token embedding space and prepending it as a global context token, the Transformer attention mechanism allows every discrete decoding step to attend to the global semantic plan. The learning objective is the negative log-likelihood over the masked regions \mathcal{M} :

$$\mathcal{L}_{\text{discrete}} = \mathbb{E}_{x,t,r} \left[- \sum_{i \in \mathcal{M}} \log p_\theta(x_i | \tilde{x}_{\setminus \mathcal{M}}, r) \right]. \quad (5)$$

This formulation fundamentally differs from RCG, which relies on pixel-space diffusion or latent VAE decoders. By operating in discrete token space, it natively handles text and quantized images (via VQ-VAE tokens) in a unified architecture.

3.4 Cross-Modal Alignment via Inter-Modal Optimal Transport

A critical challenge in multimodal generation is the misalignment between visual and textual representation spaces. We address this via a **Mixed-Modal Sampling Strategy** that effectively approximates an optimal transport plan between modalities.

We construct training batches $\mathcal{B} = \{\mathcal{B}_{\text{txt}}, \mathcal{B}_{\text{img}}, \mathcal{B}_{\text{pair}}\}$.

- **Intra-Modal Learning:** For \mathcal{B}_{txt} and \mathcal{B}_{img} , we train the model to reconstruct x given its own representation $r = \mathcal{E}(x)$.
- **Cross-Modal Bridge:** For paired data $\mathcal{B}_{\text{pair}}$, we perform **representation swapping**. We train the model to generate image tokens x_{img} conditioned on text representations r_{txt} , and vice-versa.

To facilitate this, we introduce lightweight **Modality Adapters** $\mathcal{A}_{T \rightarrow V}$ and $\mathcal{A}_{V \rightarrow T}$ that project representations into a shared semantic centroid before

injection. This forces the latent diffusion model (Stage I) and the discrete generator (Stage II) to agree on a unified semantic coordinate system, enabling zero-shot generation (e.g., text-to-image) even with limited paired data.

4 Experiments

We empirically evaluate the effectiveness of **CoM-DAD** on both unimodal and cross-modal generative tasks. Our goals are threefold: (1) to assess the generative quality and efficiency of the discrete absorbing process on high-dimensional image and text manifolds, (2) to validate the hierarchical decoupling hypothesis by analyzing the impact of the Continuous Latent Planner and Semantic Injection Interface on convergence and stability, and (3) to investigate the efficacy of Stochastic Mixed-Modal Transport for zero-shot cross-modal alignment and generalization.

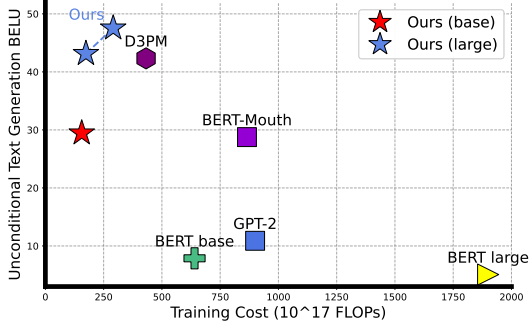
4.1 Implementation Details

Data Sources and Mixed-Modal Transport. Following our **Stochastic Mixed-Modal Transport** strategy (Sec. 3.4), we construct a unified training distribution comprising three subsets: (i) large-scale textual data from BookCorpus and Wikipedia (Wettig et al., 2022) (28M samples), (ii) image-only data from ImageNet-1k (1.28M samples), and (iii) 100K image-text pairs curated from COCO, with synthetic captions generated using Chameleon (Team, 2024). To ensure balanced manifold coverage and stable prior swapping, the sampling ratio across these sources is fixed at 2:2:1.

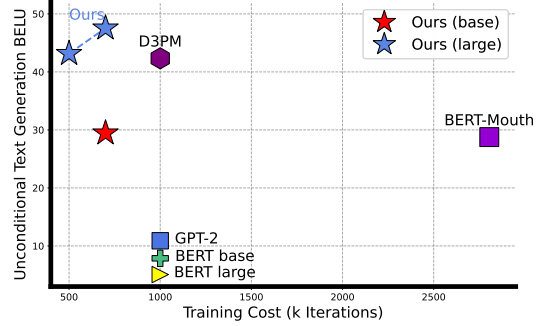
Discrete Manifold Tokenization. To establish the discrete state space, images are resized and center-cropped to 256×256 pixels, then tokenized using VQGAN (Yu et al., 2021) into 256 discrete visual tokens. Text inputs are tokenized using the RoBERTa tokenizer (Liu, 2019) and padded or truncated to a maximum of 256 tokens. No additional data augmentation is applied, relying on the variable-rate noise schedule for robustness.

Continuous Manifold and Optimization. We utilize frozen self-supervised encoders to define the continuous semantic manifold: MoCoV3 (Chen et al., 2021) for images and MPNet for text. The framework is trained in two decoupled stages:

- **Stage I (Continuous Latent Diffusion):** The continuous planner is trained 400K iterations to model the density of semantic embeddings.



(a) Training FLOPs versus BLEU performance.



(b) Convergence behavior across training iterations.

Figure 3: **Impact of the Semantic Injection Interface on convergence efficiency.** CoM-DAD achieves superior BLEU scores with substantially reduced training costs compared to the ablated variant without the interface. The reported cost includes training of both the *Continuous Latent Planner* (Stage I) and *Discrete Absorbing Diffusion* (Stage II). These results indicate that the **Semantic Injection Interface** accelerates convergence by enabling the discrete model to exploit the structure of the continuous semantic manifold, rather than learning it from scratch.

370 • **Stage II (Discrete Absorbing Diffusion):**

371 The discrete generator is trained 300K iterations
 372 to generate tokens from absorbing states.

373 We use the AdamW optimizer with an initial learn-
 374 ing rate of 5×10^{-4} . All experiments are conducted
 375 on 8 NVIDIA A800 GPUs.

376 **Evaluation Protocol.** We evaluate the topologi-
 377 cal unification capabilities of CoM-DAD. For un-
 378 conditional image generation, we follow (Li et al.,
 379 2024) and generate 50K samples from the Image-
 380 Net distribution, reporting Inception Score (IS)
 381 (Salimans et al., 2016) and Fréchet Inception Dis-
 382 tance (FID) (Heusel et al., 2017). For text gener-
 383 ation, we report BLEU (n-gram accuracy) (Papineni
 384 et al., 2002) and Self-BLEU (diversity) (Zhu et al.,
 385 2018). For cross-modal generation, we test the
 386 Modality Adapters by prompting with text to syn-
 387 thesize aligned images, evaluating semantic consis-
 388 tency through qualitative visual inspection.

389 **4.2 Main Results and Analysis**

390 **Superior fidelity on discrete text manifolds.**

391 We first evaluate CoM-DAD on unconditional
 392 text generation and compare it with strong base-
 393 lines, including autoregressive models and stan-
 394 dard masked language models. Table 1 shows
 395 that CoM-DAD achieves BLEU-2 and BLEU-4
 396 scores of 47.46 and 13.64, respectively, outper-
 397 forming all prior approaches under comparable set-
 398 tings. We also include a variant labeled “Ours
 399 (large) + Autoregressive”, where CoM-DAD op-
 400 erates in a sequential manner, demonstrating the
 401 framework’s flexibility to encompass autoregres-
 402 sive generation as a special case. Compared to stan-
 403 dard autoregressive models such as GPT-2, CoM-
 404 DAD demonstrates exceptional stability in generat-

ing long-range dependencies from scratch, validat-
 ing the efficacy of decoupling semantic planning
 from token generation. The generated text exhibits
 not only local fluency but also superior global co-
 herence, a direct result of the Continuous Latent
 Planner governing the generation trajectory.

411 **Semantic Injection Interface facilitates conver-**
 412 **gence.**

413 Figure 3 analyzes the relationship be-
 414 tween generation quality and training cost. CoM-
 415 DAD achieves superior convergence rates with sig-
 416 nificantly fewer iterations than an ablated variant
 417 lacking the **Semantic Injection Interface**. This
 418 indicates that conditioning the discrete absorbing
 419 process on the continuous semantic manifold al-
 420 lows the model to bypass the difficulty of learning
 421 structure from scratch. Furthermore, models utiliz-
 422 ing this interface produce longer and semantically
 423 richer sequences, whereas unguided counterparts
 424 tend to suffer from mode collapse or repetition,
 425 failing to bridge the gap between the discrete and
 426 continuous states effectively.

426 **Parallel decoding via Discrete Absorbing Dif-**
 427 **fusion.**

428 Unlike autoregressive models that are
 429 bound by serial token decoding, CoM-DAD lever-
 430 ages a Discrete Absorbing Diffusion process that
 431 allows for non-autoregressive parallel generation.
 432 As shown in Figure 4(a), the model can recon-
 433 struct up to 20 tokens per step simultaneously,
 434 significantly reducing inference latency. Despite
 435 the absence of aggressive GPU-specific optimiza-
 436 tions found in mature autoregressive systems, CoM-
 437 DAD achieves a $5\times$ speedup over GPT-2. This
 438 efficiency confirms that modeling generation as an
 439 iterative denoising process from absorbing states is
 440 a viable and high-throughput alternative to standard
 441 causal modeling.

Table 2: **Quantitative comparison of unconditional image generation.** CoM-DAD achieves competitive performance against state-of-the-art baselines on the continuous image manifold. These results validate the framework’s *Topological Unification*, demonstrating that the *Stochastic Mixed-Modal Transport* strategy effectively aligns discrete token generation with continuous visual semantics.

Unconditional Generation	params	FID (↓)	IS (↑)
BigGAN (Brock et al., 2018)	70M	38.61	24.7
IC-GAN (Casanova et al., 2021)	75M	15.6	59
ADM (Dhariwal and Nichol, 2021)	554M	26.21	39.7
ADDP (Tian et al., 2023)	176M	8.9	95.3
MaskGIT (Chang et al., 2022)	227M	20.72	42.1
RDM-IN (Blattmann et al., 2022)	400M	5.91	158.8
MAGE-B (Li et al., 2023)	176M	8.67	94.8
MAGE-L (Li et al., 2023)	439M	7.04	123.5
RCG-B (Li et al., 2024)	239M	3.98	177.8
RCG-L (Li et al., 2024)	502M	3.44	186.9
Ours (base)	318M	5.14	138.3
Ours (large)	593M	4.32	151.6

and the diffusion noise schedules.

Impact of the Continuous Latent Planner. To validate the necessity of topological decoupling, we remove the *Continuous Latent Planner* (Stage I) and train the *Discrete Absorbing Diffusion* model directly on token sequences without the *Semantic Injection Interface*. As illustrated in Figure 6a, models trained without continuous latent conditioning can produce short, syntactically correct phrases but often fail when tasked with generating longer or semantically complex passages. Furthermore, they require significantly more training iterations to achieve comparable quality. This supports our fundamental hypothesis that externalizing semantic planning into a continuous manifold simplifies the discrete learning objective, allowing the token generator to focus solely on mapping high-level plans to discrete structures.

Necessity of the Variable-Rate Noise Schedule. We retrain CoM-DAD using a standard fixed 15% masking rate (typical of BERT-style MLMs). As shown in Figure 6b, this variant fails to generate coherent text from the fully absorbing state, instead repeating simple or semantically trivial tokens. In contrast, the **Variable-Rate Noise Schedule** employed in CoM-DAD successfully generates complete sentences from scratch. This confirms that aggressive, variable-rate masking is essential for shifting the model from a local infilling objective to a global generative task. These findings are consistent with our insight in Sec. 3.3 and provide empirical support for *Discrete Absorbing Diffusion* as a mechanism for robust generation rather than mere masked prediction.

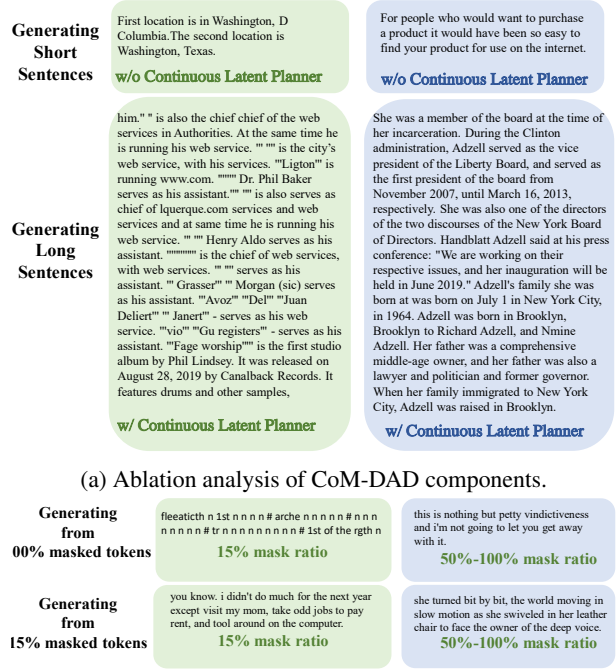


Figure 6: **Ablation studies of CoM-DAD.** (a) Removing the **Continuous Latent Planner** (and the corresponding **Semantic Injection Interface**) substantially degrades the model’s ability to generate complex, long-form passages, leading to slower convergence and reduced global coherence. This supports our hypothesis that topological decoupling simplifies the learning objective for the **Discrete Absorbing Diffusion** process. (b) Comparison of the **Variable-Rate Noise Schedule** against a fixed 15% rate. The fixed-rate variant yields repetitive, trivial outputs when starting from a *fully absorbing state*. In contrast, CoM-DAD’s variable schedule enables generation from scratch, proving that dynamic masking is critical to shift from local infilling to global generation.

5 Conclusion

Summary of Benefits. CoM-DAD’s architectural and training design confers several benefits: (1) **Efficiency:** By decoupling representation modeling from token generation, CoM-DAD enables faster convergence and 5× faster inference over standard denoising or autoregressive models. (2) **Generative Capability:** The variable masking schedule fosters the model’s ability to generate coherent and diverse outputs rather than merely recover corrupted inputs. (3) **Multimodal Alignment:** Our mixed sampling strategy facilitates scalable training from unimodal data while achieving strong cross-modal consistency. (4) **Unified Architecture:** A single encoder-decoder model handles both text and image generation through a shared conditioning mechanism, supporting flexible and generalizable generation tasks.

534 Limitations

535 While CoM-DAD effectively bridges the topologi-
536 cal gap between discrete and continuous modalities,
537 our current empirical validation is primarily fo-
538 cused on foundational image-text generation tasks.
539 Although the Stochastic Mixed-Modal Transport
540 framework is theoretically extensible to temporal
541 modalities like video or audio, we reserve the spe-
542 cific calibration of the Variable-Rate Noise Sched-
543 ule for these high-dimensional domains for future
544 work to maintain focused analysis. Furthermore,
545 we observe that standard automated metrics may
546 not fully capture the long-horizon semantic con-
547 sistency driven by the Continuous Latent Planner,
548 potentially underrepresenting the model’s ability
549 to generate conceptually accurate but structurally
550 diverse outputs compared to rigid token-matching
551 baselines.

552 References

553 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc,
554 Antoine Miech, Iain Barr, Yana Hasson, Karel
555 Lenc, Arthur Mensch, Katherine Millican, Malcolm
556 Reynolds, and 1 others. 2022. Flamingo: a visual
557 language model for few-shot learning. *Advances in*
558 *neural information processing systems*, 35:23716–
559 23736.

560 Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel
561 Tarlow, and Rianne Van Den Berg. 2021. Structured
562 denoising diffusion models in discrete state-spaces.
563 *Advances in Neural Information Processing Systems*,
564 34:17981–17993.

565 Andreas Blattmann, Robin Rombach, Kaan Oktay,
566 Jonas Müller, and Björn Ommer. 2022. Retrieval-
567 augmented diffusion models. *Advances in Neural*
568 *Information Processing Systems*, 35:15309–15324.

569 Andrew Brock, Jeff Donahue, and Karen Simonyan.
570 2018. Large scale gan training for high fi-
571 delity natural image synthesis. *arXiv preprint*
572 *arXiv:1809.11096*.

573 Arantxa Casanova, Marlene Careil, Jakob Verbeek,
574 Michal Drozdal, and Adriana Romero Soriano.
575 2021. Instance-conditioned gan. *Advances in Neural*
576 *Information Processing Systems*, 34:27517–27529.

577 Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and
578 William T Freeman. 2022. Maskgit: Masked gener-
579 ative image transformer. In *Proceedings of the*
580 *IEEE/CVF Conference on Computer Vision and Pat-*
581 *tern Recognition*, pages 11315–11325.

582 Xi Chen, Josip Djolonga, Piotr Padlewski, Basil
583 Mustafa, Soravit Changpinyo, Jialin Wu, Car-
584 los Riquelme Ruiz, Sebastian Goodman, Xiao Wang,

585 Yi Tay, and 1 others. 2023. Pali-x: On scaling up
586 a multilingual vision and language model. *arXiv*
587 *preprint arXiv:2305.18565*.

Xinlei Chen, Saining Xie, and Kaiming He. 2021. An
588 empirical study of training self-supervised vision
589 transformers. In *Proceedings of the IEEE/CVF in-*
590 *ternational conference on computer vision*, pages
591 9640–9649. 592

Jacob Devlin. 2018. Bert: Pre-training of deep bidi-
593 rectional transformers for language understanding.
594 *arXiv preprint arXiv:1810.04805*. 595

Prafulla Dhariwal and Alexander Nichol. 2021. Diffu-
596 sion models beat gans on image synthesis. *Advances*
597 *in neural information processing systems*, 34:8780–
598 8794. 599

Google DeepMind. 2025. Gemini diffusion. <https://deepmind.google/models/gemini-diffusion>.
600 Accessed: 2025-05-24. 601 602

Zhengfu He, Tianxiang Sun, Kuanning Wang, Xuan-
603 jing Huang, and Xipeng Qiu. 2022. Diffusionbert:
604 Improving generative masked language models with
605 diffusion models. *arXiv preprint arXiv:2211.15029*. 606

Martin Heusel, Hubert Ramsauer, Thomas Unterthiner,
607 Bernhard Nessler, and Sepp Hochreiter. 2017. Gans
608 trained by a two time-scale update rule converge to a
609 local nash equilibrium. *Advances in neural informa-*
610 *tion processing systems*, 30. 611

Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. De-
612 noising diffusion probabilistic models. *Advances*
613 *in neural information processing systems*, 33:6840–
614 6851. 615

Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini,
616 Patrick Forré, and Max Welling. 2021. Argmax flows
617 and multinomial diffusion: Learning categorical dis-
618 tributions. *Advances in Neural Information Process-*
619 *ing Systems*, 34:12454–12465. 620

Tianhong Li, Huiwen Chang, Shlok Mishra, Han Zhang,
621 Dina Katabi, and Dilip Krishnan. 2023. Mage:
622 Masked generative encoder to unify representation
623 learning and image synthesis. In *Proceedings of the*
624 *IEEE/CVF Conference on Computer Vision and Pat-*
625 *tern Recognition*, pages 2142–2152. 626

Tianhong Li, Dina Katabi, and Kaiming He. 2024. Re-
627 turn of unconditional generation: A self-supervised
628 representation generation method. In *The Thirty-*
629 *eighth Annual Conference on Neural Information*
630 *Processing Systems*. 631

Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S
632 Liang, and Tatsunori B Hashimoto. 2022. Diffusion-
633 lm improves controllable text generation. *Advances*
634 *in Neural Information Processing Systems*, 35:4328–
635 4343. 636

637	Xiao Li, Jiaqi Zhang, Shuxiang Zhang, Tianshui Chen,	Chengyue Wu, Hao Zhang, Shuchen Xue, Zhijian Liu,	689
638	Liang Lin, and Guangrun Wang. 2025. In-situ	Shizhe Diao, Ligeng Zhu, Ping Luo, Song Han, and	690
639	tweedie discrete diffusion models. <i>arXiv preprint</i>	Enze Xie. 2025. Fast-dllm: Training-free accelera-	691
640	<i>arXiv:2510.01047</i> .	tion of diffusion llm by enabling kv cache and parallel	692
		decoding. <i>arXiv preprint arXiv:2505.22618</i> .	693
641	Yinhan Liu. 2019. Roberta: A robustly opti-	Jiacheng Ye, Zihui Xie, Lin Zheng, Jiahui Gao, Zirui	694
642	mized bert pretraining approach. <i>arXiv preprint</i>	Wu, Xin Jiang, Zhenguo Li, and Lingpeng Kong.	695
643	<i>arXiv:1907.11692</i> , 364.	2025. Dream 7b .	696
644	Shen Nie, Fengqi Zhu, Zebin You, Xiaolu Zhang,	Jiahui Yu, Xin Li, Jing Yu Koh, Han Zhang, Ruom-	697
645	Jingyang Ou, Jun Hu, Jun Zhou, Yankai Lin, Ji-Rong	ing Pang, James Qin, Alexander Ku, Yuanzhong Xu,	698
646	Wen, and Chongxuan Li. 2025. Large language dif-	Jason Baldrige, and Yonghui Wu. 2021. Vector-	699
647	fusion models. <i>arXiv preprint arXiv:2502.09992</i> .	quantized image modeling with improved vqgan.	700
		<i>arXiv preprint arXiv:2110.04627</i> .	701
648	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan	702
649	Jing Zhu. 2002. Bleu: a method for automatic evalua-	Zhang, Jun Wang, and Yong Yu. 2018. Texus: A	703
650	tion of machine translation. In <i>Proceedings of the</i>	benchmarking platform for text generation models.	704
651	<i>40th annual meeting of the Association for Computa-</i>	In <i>The 41st international ACM SIGIR conference</i>	705
652	<i>tional Linguistics</i> , pages 311–318.	<i>on research & development in information retrieval</i> ,	706
		pages 1097–1100.	707
653	Tim Salimans, Ian Goodfellow, Wojciech Zaremba,		
654	Vicki Cheung, Alec Radford, and Xi Chen. 2016.		
655	Improved techniques for training gans. <i>Advances in</i>		
656	<i>neural information processing systems</i> , 29.		
657	Jascha Sohl-Dickstein, Eric Weiss, Niru Mah-		
658	eswaranathan, and Surya Ganguli. 2015. Deep un-		
659	supervised learning using nonequilibrium thermo-		
660	dynamics. In <i>International conference on machine</i>		
661	<i>learning</i> , pages 2256–2265. PMLR.		
662	Chameleon Team. 2024. Chameleon: Mixed-modal		
663	early-fusion foundation models. <i>arXiv preprint</i>		
664	<i>arXiv:2405.09818</i> .		
665	Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-		
666	Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan		
667	Schalkwyk, Andrew M Dai, Anja Hauth, Katie Mil-		
668	lican, and 1 others. 2023. Gemini: a family of		
669	highly capable multimodal models. <i>arXiv preprint</i>		
670	<i>arXiv:2312.11805</i> .		
671	Changyao Tian, Chenxin Tao, Jifeng Dai, Hao Li, Zi-		
672	heng Li, Lewei Lu, Xiaogang Wang, Hongsheng Li,		
673	Gao Huang, and Xizhou Zhu. 2023. Addp: Learn-		
674	ing general representations for image recognition and		
675	generation with alternating denoising diffusion pro-		
676	cess. <i>arXiv preprint arXiv:2306.05423</i> .		
677	Alex Wang and Kyunghyun Cho. 2019. Bert has		
678	a mouth, and it must speak: Bert as a markov		
679	random field language model. <i>arXiv preprint</i>		
680	<i>arXiv:1902.04094</i> .		
681	Guangrun Wang and Philip HS Torr. 2022. Traditional		
682	classification neural networks are good generators:		
683	They are competitive with ddpms and gans. <i>arXiv</i>		
684	<i>preprint arXiv:2211.14794</i> .		
685	Alexander Wettig, Tianyu Gao, Zexuan Zhong, and		
686	Danqi Chen. 2022. Should you mask 15% in		
687	masked language modeling? <i>arXiv preprint</i>		
688	<i>arXiv:2202.08005</i> .		