

Diffuse Everything : Multimodal Diffusion Models on Arbitrary State Spaces

Kevin Rojas*¹ Yuchen Zhu*¹ Sichen Zhu¹ Felix X.-F. Ye² Molei Tao¹

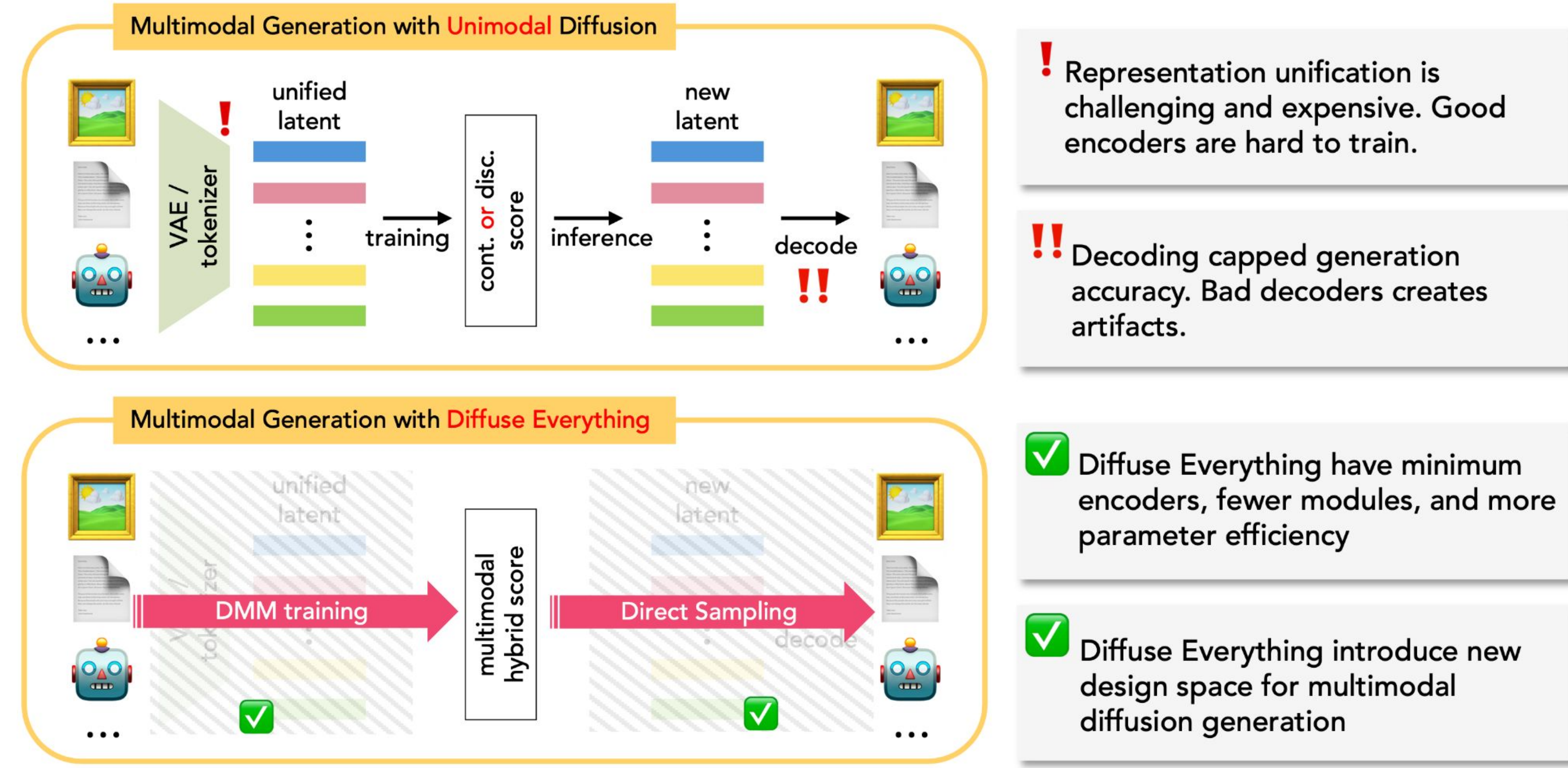
¹Georgia Institute of Technology ²SUNY Albany



International Conference
On Machine Learning



TL;DR



Motivation

- Joint generation of multimodal data is **IMPORTANT!**
- Diffusion models are **SOTA** for many types of unimodal data, with fantastic conditional generation capability.
- Good tokenizers and VAEs are **CHALLENGING** to train.

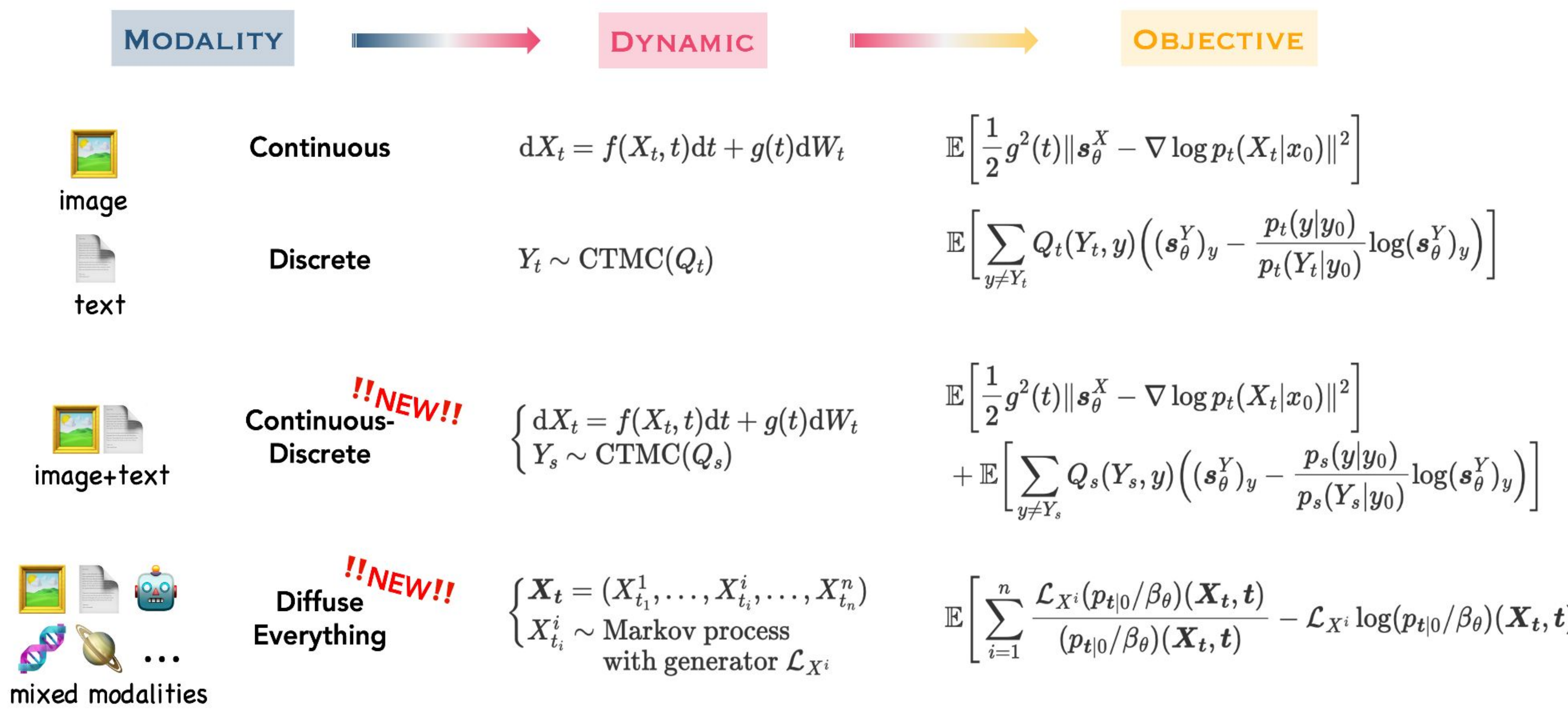
Goal

Develop **NEW** diffusion models to generate multimodal data in their **native state space**, **bypassing** the need for tokenizers/VAEs/encoders.

Key Findings

- Multimodal diffusions can be built by **combining** unimodal diffusions and be trained by learning scores of **joint distribution**.
- Training multimodal diffusion models is **provably** as simple as jointly optimizing **sum** of unimodal diffusion model losses.
- Decoupled time enables **any-to-any** generation in one model, and a new guidance scheme named noisy guidance.
- Multimodal diffusion on native state spaces are much more **parameter-efficient**.

Approach



Methodology

Multimodal Gen. with Denoising Markov Models

Generative modeling with dynamic:

$$\mathbf{X}_t = (X_t^1, \dots, X_t^i, \dots, X_t^n), 0 \leq t^1, \dots, t^n \leq T$$
$$(X_0^1, \dots, X_0^i, \dots, X_0^n) \sim p_{\text{data}}(\mathbf{x})$$

... requires only learning score of **Theorem 1**

$$p(\mathbf{x}_t, t) = \mathbb{P}(X_t^1 = (\mathbf{x}_t)_1, \dots, X_t^n = (\mathbf{x}_t)_n)$$

Training with Generalized Score Matching

$$\mathcal{I}_{\text{GESM}} = \mathbb{E}_{t, \mathbf{x}_t \sim p(\cdot, t)} \left[\sum_{i=0}^n \frac{\mathcal{L}_{X^i}(p/\beta_\theta)(\mathbf{x}_t, t)}{(p/\beta_\theta)(\mathbf{x}_t, t)} - \mathcal{L}_{X^i} \log(p/\beta_\theta)(\mathbf{x}_t, t) \right]$$
$$\mathcal{I}_{\text{GDMS}} = \mathbb{E}_{t, p_0, p_{t|0}} \left[\sum_{i=1}^n \frac{\mathcal{L}_{X^i}(p_{t|0}/\beta_\theta)(\mathbf{x}_t, t)}{(p_{t|0}/\beta_\theta)(\mathbf{x}_t, t)} - \mathcal{L}_{X^i} \log(p_{t|0}/\beta_\theta)(\mathbf{x}_t, t) \right]$$
$$\mathcal{I}_{\text{GISM}} = \mathbb{E}_{t, p_t} \left[\sum_{i=1}^n \frac{\mathcal{L}_{X^i}^*(\beta_\theta)(\mathbf{x}_t, t)}{\beta_\theta(\mathbf{x}_t, t)} - \mathcal{L}_{X^i}^* \log(\beta_\theta)(\mathbf{x}_t, t) \right]$$

Equivalent
Theorem 2

Continuous-Discrete Diffusion Model

Forward Process

$$\begin{cases} d\tilde{X}_t = f(\tilde{X}_t, t)dt + g(t)dB_t \\ Y_s \sim \text{CTMC}(Q_s), (\tilde{X}_0, Y_0) \sim p_{\text{data}}(x, y) \end{cases}$$

Backward Process

$$\begin{cases} d\tilde{X}_t = -f(\tilde{X}_t, T-t) + g^2(T-t) \nabla \log \tilde{p}_{t,s}(\tilde{X}_t, \tilde{Y}_s) dt + g(T-t)dB_t \\ \tilde{Y}_s \sim \text{CTMC}(\tilde{Q}(\tilde{X}_t, t, s)), (\tilde{X}_0, \tilde{Y}_0) \sim p(x, y, T, T) \\ \tilde{p}_{t,s} = p_{T-t, T-s}, \tilde{Q}(\tilde{X}_t, t, s)_{y,y'} = \frac{\tilde{p}_{t,s}(\tilde{X}_t, y)}{\tilde{p}_{t,s}(\tilde{X}_t, y')} (Q_{T-t})_{y,y'} \end{cases}$$

Score of joint dist.

- A combination of diffusion SDE on Euclidean space and CTMC jump process on finite state space.
- Aiming at generating data with both continuous values and discrete values

Denoising training objective

$$\mathbb{E}_{t, s, x_0, y_0 \sim p_0, x_t, y_t \sim p_{t|0}} \left[\frac{1}{2} g^2(t) \|s_\theta^X - \nabla \log p_t(x_t|x_0)\|^2 + \sum_{y \neq y_s} Q_s(y, y) \left((s_\theta^Y)_y - \frac{p_s(y|y_0)}{p_s(y_s|y_0)} \log(s_\theta^Y)_y \right) \right]$$

Requires multimodal input
Unimodal conditional score

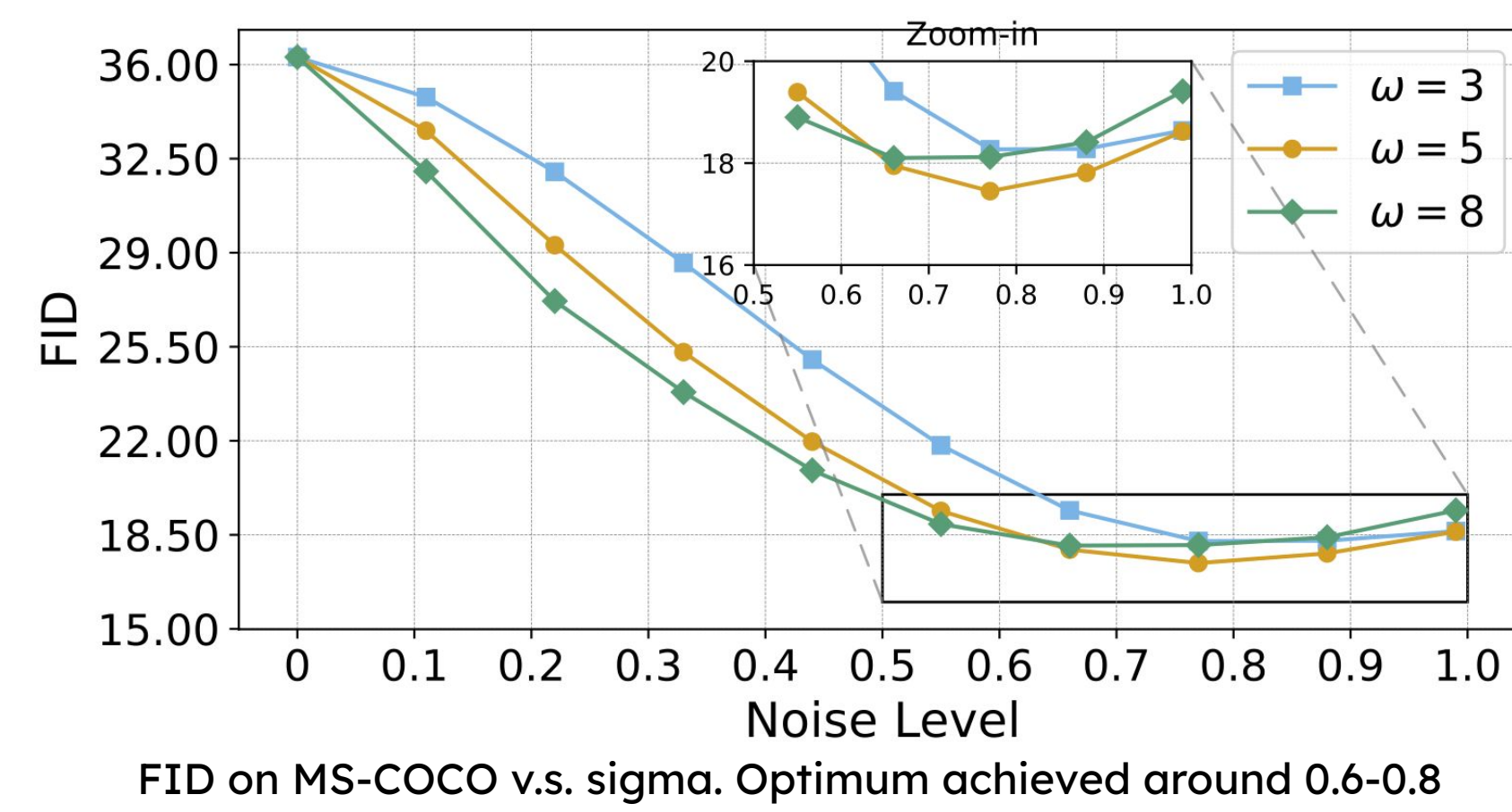
Decoupled Time Design

Decoupled time design means:

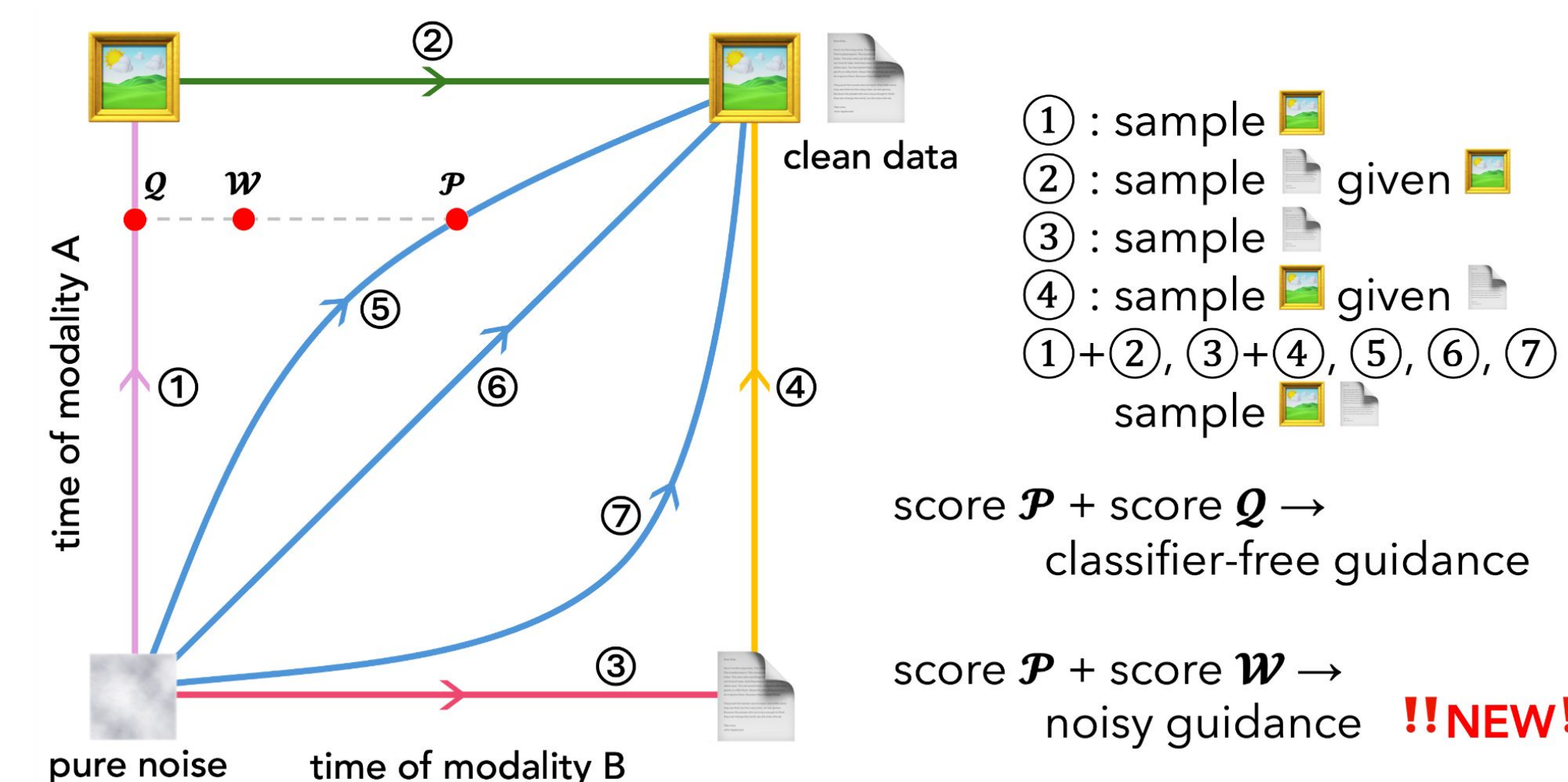
- Each modality is noised and denoised at independent pace
- Requires learning denoiser/score at more scenarios

Benefit 1 - Flexible any-to-any generation

- Given a partially noisy text Y_s , simulating only the X-backward dynamics samples from $p_{\text{data}}(x|Y_s, s)$
- Given a partially noisy image X_t , simulating the Y-backward dynamics samples from $p_{\text{data}}(y|X_t, t)$



FID on MS-COCO v.s. sigma. Optimum achieved around 0.6-0.8



Benefit 2 - Better guidance scheme than CFG

- Noisy guidance: guiding with a **partial corrupted conditional** model
- Achieving better diversity-quality trade-off
- Recover CFG when $s = 0$, $\sigma = T$

Noisy Guidance

$$\omega s_\theta(x_t, y_s, t, s) + (1 - \omega) s_\theta(x_t, y_\sigma, t, \sigma) \quad \sigma > s$$

Inference Algorithms

Algorithm 3 Discrete Sampler with τ -leaping

Require: N : Number of steps, ω : Guidance Strength

1: $[a, b]$: Guidance Interval, model: $s_\theta(x_t, y_s, t, s, \omega)$

2: y : A clean image condition

Ensure: $y_0 \sim p_{\text{data}}(\cdot|x)$

3: $y_t \leftarrow [M, \dots, M]$

4: **for** t in times **do**

5: $\omega_t = \omega$ if $t \in [a, b]$ else 1.

6: $s_\theta^x, s_\theta^y \leftarrow s_\theta(x_t, y_t, t, \omega_t)$

7: $v_{\text{old}} = f(x_t, t) - \frac{1}{2} g^2(t) s_\theta^x$

8: $\tilde{x} \leftarrow x_t + v_{\text{old}} dt, \tilde{s}_\theta^x, \tilde{s}_\theta^y \leftarrow s_\theta(\tilde{x}, y_t, t + dt, \omega_t)$

9: $v_{\text{new}} = f(\tilde{x}, t) - \frac{1}{2} g^2(t) \tilde{s}_\theta^x$

10: $x_t \leftarrow x_t + \frac{1}{2} \cdot (v_{\text{old}} + v_{\text{new}}) dt$

11: **end for**

12: **return** x_0

Algorithm 4 Continuous Sampler with Heun's method

Require: N : Number of steps, ω : Guidance Strength

1: $[a, b]$: Guidance Interval, model: $s_\theta(x_t, y_s, t, s, \omega)$

2: model: $s_\theta(x_t, y_s, t, s, \omega)$

Ensure: $x_0, y_0 \sim p_{\text{data}}(\cdot|x)$

3: $x_t \leftarrow \mathcal{N}(0, I), y_t \leftarrow [M, \dots, M]$

4: **for** t in times **do**

5: $\omega_t = \omega$ if $t \in [a, b]$ else 1.

6: $s_\theta^x, s_\theta^y \leftarrow s_\theta(x_t, y_t, t, \omega_t), v_{\text{old}} = f(x_t, t) - \frac{1}{2} g^2(t) s_\theta^x$

7: $\tilde{x} \leftarrow x_t + v_{\text{old}} dt$

8: $s_\theta^x, s_\theta^y \leftarrow s_\theta(\tilde{x}, y_t, t + dt, \omega_t), v_{\text{new}} = f(\tilde{x}, t) - \frac{1}{2} g^2(t) \tilde{s}_\theta^x$

9: $x_t \leftarrow x_t + \frac{1}{2} \cdot (v_{\text{old}} + v_{\text{new}}) dt$

10: **end for**

11: **return** x_0, y_0

Algorithm 5 Multimodal Sampler with τ -leaping and Heun's Method

Require: N : Number of steps, ω : Guidance Strength,

1: $[a, b]$: Guidance Interval,

2: model: $s_\theta(x_t, y_s, t, s, \omega)$

Ensure: $x_0, y_0 \sim p_{\text{data}}(\cdot|x)$

3: $x_t \leftarrow \mathcal{N}(0, I), y_t \leftarrow [M, \dots, M]$

4: **for** t in times **do**

5: $\omega_t = \omega$ if $t \in [a, b]$ else 1.

6: $s_\theta^x, s_\theta^y \leftarrow s_\theta(x_t, y_t, t, \omega_t), v_{\text{old}} = f(x_t, t) - \frac{1}{2} g^2(t) s_\theta^x$

7: $\tilde{x} \leftarrow x_t + v_{\text{old}} dt$

8: $s_\theta^x, s_\theta^y \leftarrow s_\theta(\tilde{x}, y_t, t + dt, \omega_t), v_{\text{new}} = f(\tilde{x}, t) - \frac{1}{2} g^2(t) \tilde{s}_\theta^x$

9: $x_t \leftarrow x_t + \frac{1}{2} \cdot (v_{\text{old}} + v_{\text{new}}) dt$

10: **end for**

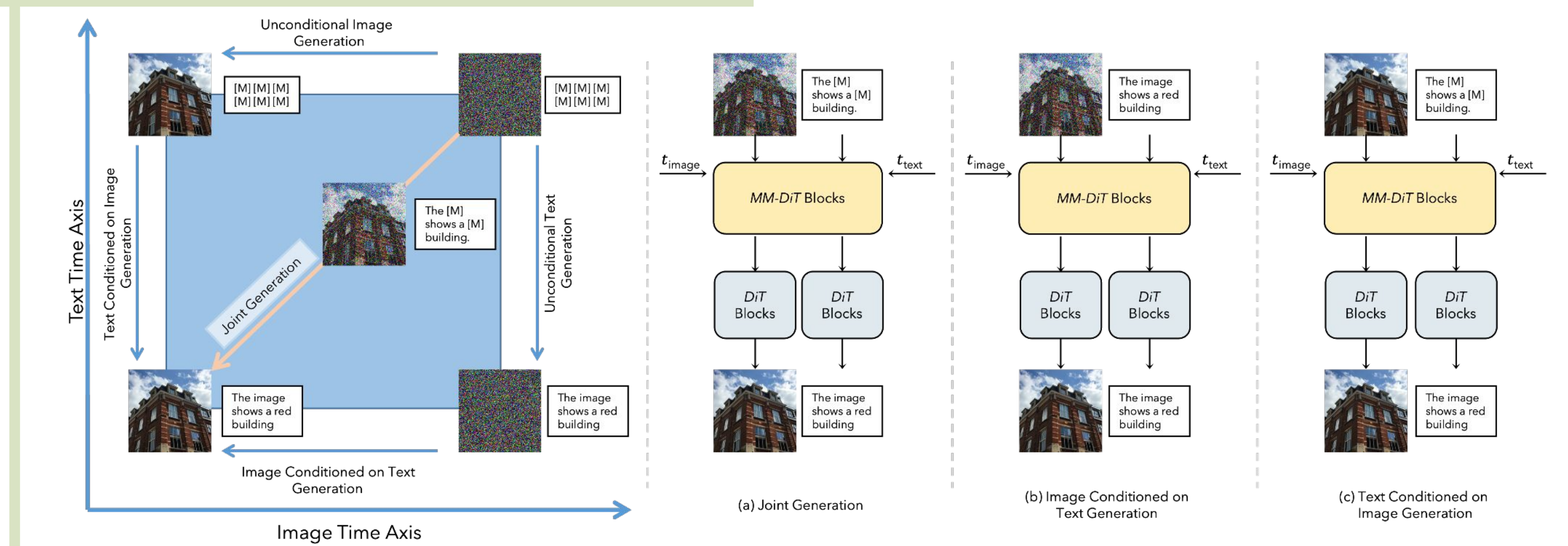
11: **return** x_0, y_0

Sampling discrete data only

Sample continuous data only

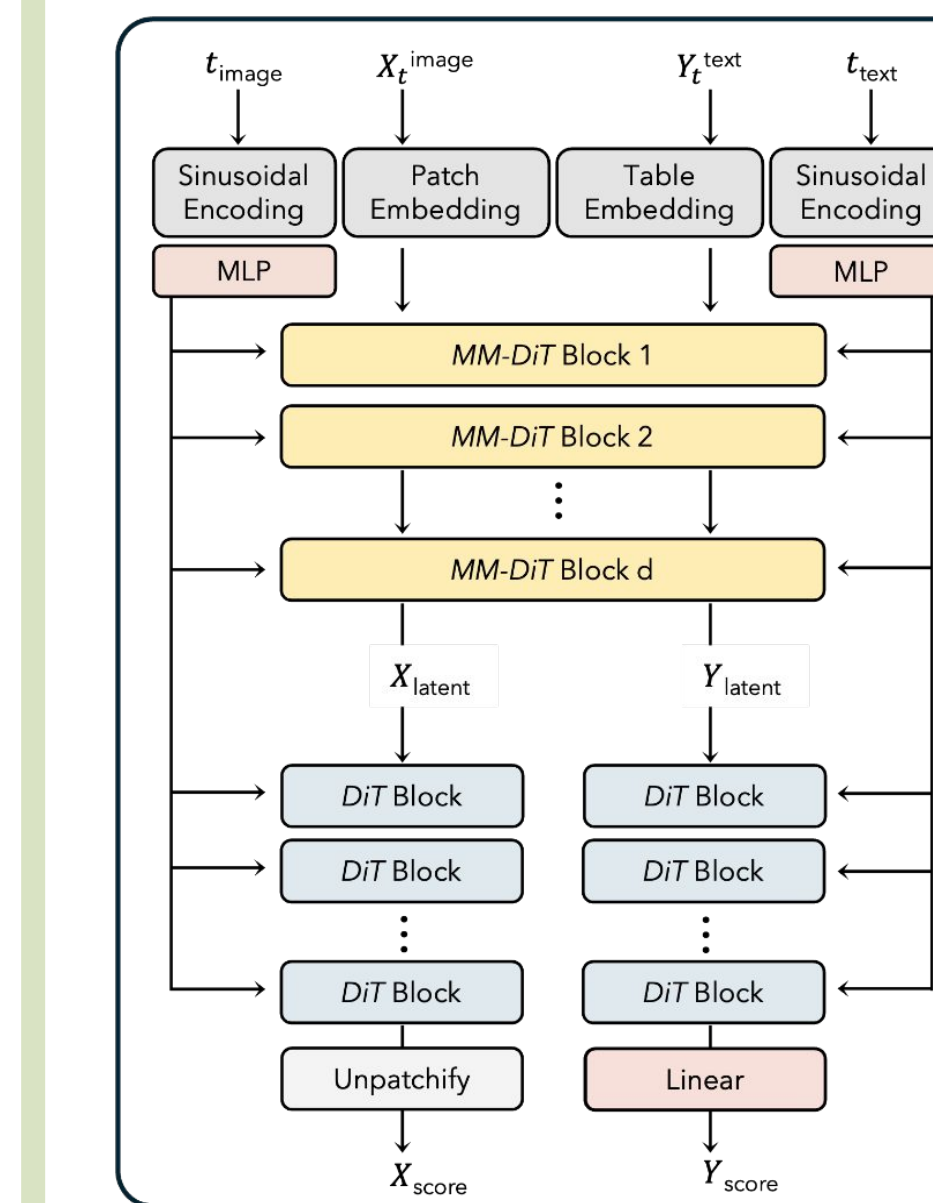
Jointly sample continuous and discrete data

Text-Image generation



Jointly generate images and its captions, and we..

- Minimally rely on pretrained model** (except for one image VAE for dimension reduction).
- No T5/CLIP/ViT/etc.
- Achieve satisfactory performance while using **much smaller model**
- Design **multi-stage training** strategy to aid training of decoupled time



score backbone

Table 1. Results on the **text to image conditional generation** on MS-COCO. We mark the extra encoders leveraged by each model with the corresponding sizes and types. SR: super resolution, TE: text encoder, VAE: variational autoencoder, VE: visual encoder, VQ-GAN: Vector Quantized GAN, VQ-VAE: vector-quantized variational autoencoder.

Model	FID	Number of Images	#Params	Extra Encoders
Models for Text-to-Image generation only				
DALL-E 2 (Ramesh et al., 2022)	10.39	650M	6.5B	123M (TE) + 700M (SR)
Imagen (Saharia et al., 2022)	7.27	860M	3B	4.6B (TE) + 600M (SR)
Stable Diffusion (Rombach et al., 2022)	12.63	400M	1.45B	123M (TE) + 83M (VAE)
PixArt- α XL/2 (Chen et al., 2023)	7.32	25M	600M	123M (TE) + 83M (VAE)
MMDiT-improved (Ifriqi et al., 2024)	6.79	12M	600M	123M (TE) + 83M (VAE)
Models for multimodal generation and understanding				
Show-o (Xie et al., 2024)	9.24	35M	1.3B	115M (VE) + 307M (VQ-VAE)
Transfusion (Zhou et al., 2024)	6.78	692M	7B	86M (VAE)
Chameleon (Meta, 2024)	26.7	600M	7B	307M (VQ-GAN)
JetFormer (Tschannen et al., 2024)	20.86	1B	2.75B	—
Models for multimodal generation only				
Versatile Diffusion (Xu et al., 2023)	11.10	400M	1.45B	123M (TE) + 83M (VAE) + 110M (TE)
UniD3 (Hu et al., 2022)	25.11	592K	600M	123M (TE) + 307M (VQ-GAN)
Our model	16.16	12M	481M	83M (VAE)

FID-30k evaluation on MS-COCO. Our model trained on SAM-LLaVA.

Mixed-type Tabular data Synthesis

Jointly generate tabular data with both categorical values (eg. age) and continuous values (eg. income), and we ...

- Achieve comparable or beat previous SOTA method with a **significantly smaller model**
- Design a new score backbone based on modification of DiT for mixed-type data, which is **effective!**

Table 2. Performance on the **Trend** metric in percentage (%). Higher values indicate better performance. Best performance in **bold**. Second best in underline.

Methods	#Parameters	Adult	Default	Shoppers	Magic	Beijing	News
GOGGLE (Liu et al., 2023)	~ 5.6M	54.71	78.06	76.10	90.53	54.06	76.81
STaSy (Kim et al., 2022)	~ 10.3M	85.49 \pm 0.25	94.04 \pm 0.26	91.51 \pm 0.15	93.39 \pm 0.53	92.00 \pm 0.10	96.93 \pm 0.04
CoDi (Lee et al., 2023)	~ 25.0M	77.51 \pm 0.08	31.59 \pm 0.05	82.22 \pm 0.11	93.47 \pm 0.25	92.93 \pm 0.15	88.90 \pm 0.01
TabDDPM (Kotelnikov et al., 2023)	~ 11.7M	96.99 \pm 0.25	95.11 \pm 0.10	93.39 \pm 0.16	98.30 \pm 0.22	97.20 \pm 0.09	86.84 \pm 0.11
TABSYN (Zhang et al., 2023)	~ 10.7M	98.46 \pm 0.27	97.93 \pm 0.12	97.93 \pm 0.21	98.94 \pm 0.31	97.76 \pm 0.28	98.56 \pm 0.03
TABSYN (reproduced)	~ 10.7M	98.29 \pm 0.22	95.25 \pm 0.51	97.82 \pm 0.14	99.16 \pm 0.16	94.86 \pm 0.34	98.52 \pm 0.09
Our model	~ 64K	98.75 \pm 0.09	96.00 \pm 1.23	98.24 \pm 0.13	98.85 \pm 0.42	97.42 \pm 0.11	98.57 \pm 0.16

Previous model size **10M~25M**

100-200X Reduction >>> Ours: **64K**

Table 3. Performance on the **MLE** metric. Higher values in AUC and lower values in RMSE indicate better testing performance. Best performance in **bold**. Second best in underline.

Methods	#Parameters	Adult (AUC \uparrow)	Default (AUC \uparrow)	Shoppers (AUC \uparrow)	Magic (AUC \uparrow)	Beijing (RMSE \downarrow)	News (RMSE \downarrow)
GOGGLE (Liu et al., 2023)	~ 5.6M	.778 \pm 0.012	.584 \pm 0.005	.658 \pm 0.052	.654 \pm 0.024	1.090 \pm 0.025	.877 \pm 0.002
STaSy (Kim et al., 2022)	~ 10.3M	.906 \pm 0.001	.752 \pm 0.006	.914 \pm 0.005	.934 \pm 0.003	.656 \pm 0.014	.871 \pm 0.002
CoDi (Lee et al., 2023)	~ 25.0M	.871 \pm 0.006	.525 \pm 0.006	.865 \pm 0.006	.932 \pm 0.003	.818 \pm 0.021	1.21 \pm 0.005
TabDDPM (Kotelnikov et al., 2023)	~ 11.7M	.907 \pm 0.001	.758 \pm 0.004	.918 \pm 0.005	.935 \pm 0.003	.592 \pm 0.011	4.86 \pm 3.04
TABSYN (Zhang et al., 2023)	~ 10.7M	.915 \pm 0.002	.764 \pm 0.004	.920 \pm 0.005	.938 \pm 0.002	.582 \pm 0.008	.861 \pm 0.027
TABSYN (reproduced)	~ 10.7M	.910 \pm 0.001	.755 \pm 0.004	.916 \pm 0.004	.939 \pm 0.003	.655 \pm 0.012	.851 \pm 0.024
Our model	~ 64K	.915 \pm 0.001	.764 \pm 0.002	.924 \pm 0.003	.941 \pm 0.002	.543 \pm 0.012	.864 \pm 0.021

MLE is the testing accuracy of the classification or regression task on real data after training an XGBoost Classifier or an XGBoost Regressor on the synthetic tabular data.