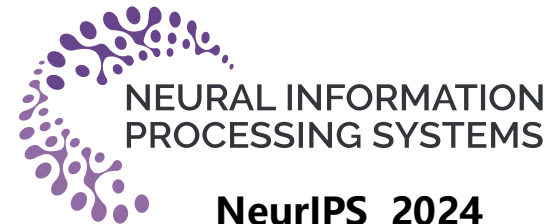




清华大学



衍远科技



# Neural Residual Diffusion Models for Deep Scalable Vision Generation

The Thirty-Eighth Annual Conference on Neural Information Processing Systems

**Zhiyuan Ma**

**(mzyth@tsinghua.edu.cn)**

Tsinghua University, Beijing, China

December, 2024

# CONTENTS



- 01 Background
- 02 Motivation
- 03 Methods
- 04 Experiment
- 05 Conclusion

# CONTENTS



## 01 Background

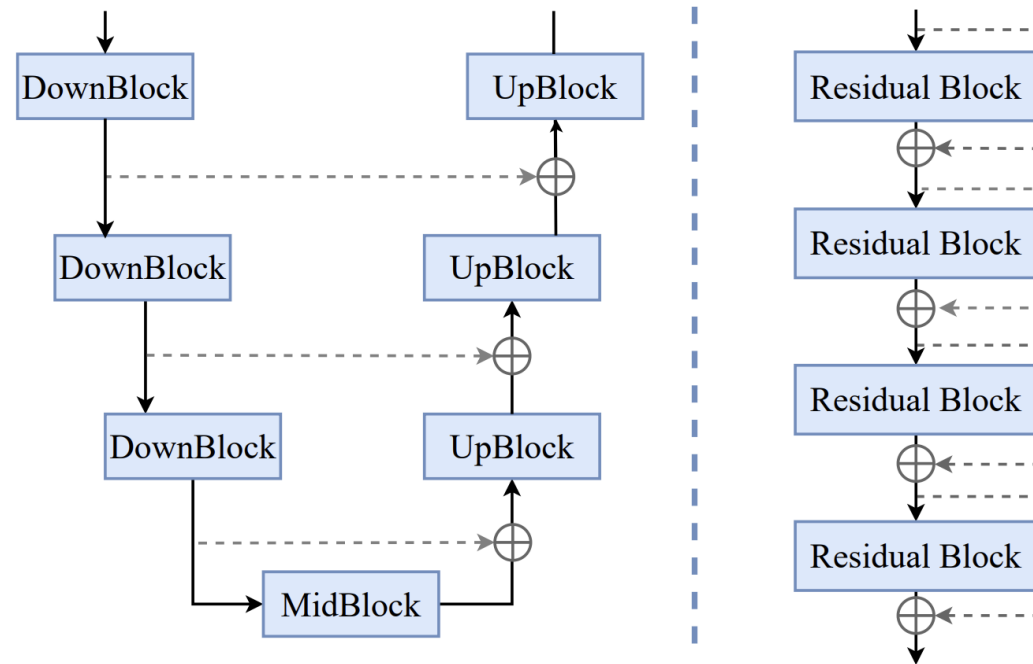
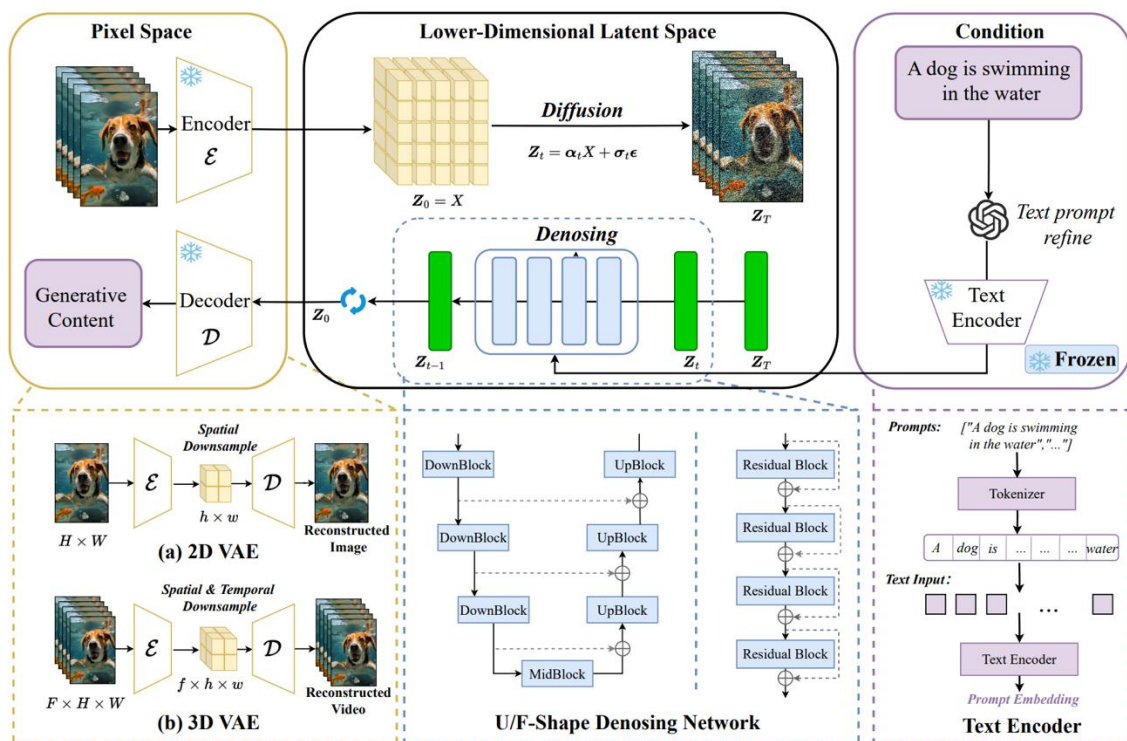
02 Motivation

03 Methods

04 Experiment

05 Conclusion

## □ Deep Generative Diffusion Networks

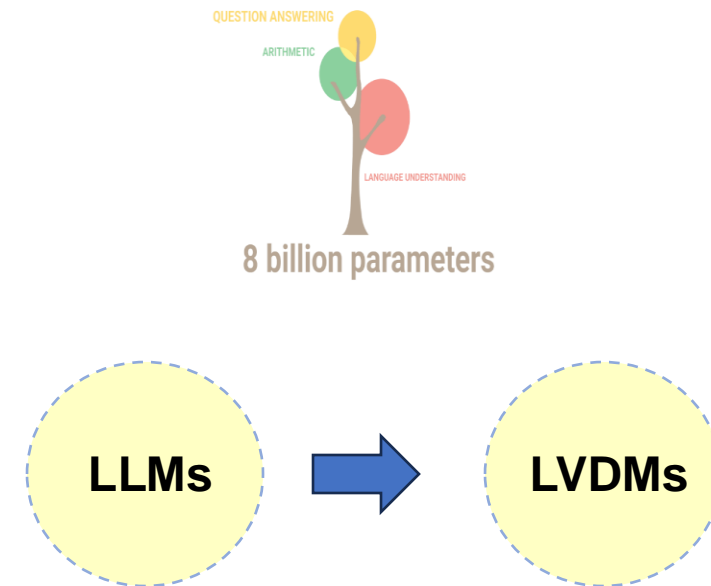


**U/F-Shape Denosing Network**

✓ The mainstream denoising backbones: U-Net, Transformer, U-ViT, DiT... (U-shaped / F-shaped)

## □ Representative Generative Diffusion Models

Methods	Year	Organization	Backbone	VAE	Text Encoder	# Params	
ADM [4]	2021	OpenAI	Unet	None	-	554M	
CDM [157]	2021	Google			-	-	
DALL-E 2 [6]	2022	OpenAI			CLIP	6.5B	
Imagen [5]	2022	Google			T5-XXL	3B	
LDM [33]	2022	LMU Munich			CLIP ViT-L	400M+55M(VAE)	
SD1.5 [33]	2022	LMU Munich	2D VAE	CLIP ViT-L	860M		
SD2.0 [33]	2022	LMU Munich		OpenCLIP ViT-H	865M		
SDXL [8]	2023	Stability AI		CLIP ViT-L & OpenCLIP ViT-bigG	2.6B		
Playground-v2.5 [158]	2024	Playground		CLIP	-		
UViT [72]	2022	Tsinghua University	Transformer		CLIP ViT-L	501M+84M(VAE)	
DiT [73]	2022	UC Berkeley			CLIP ViT-L	675M+84M(VAE)	
PixArt- $\alpha$ [81]	2023	Huawei Noah's Ark Lab			T5-XXL	600M	
FiT [74]	2024	Shanghai AI Lab			CLIP ViT-L	-	
SiT [75]	2024	New York University			CLIP ViT-L	675M	
Latte [79]	2024	Shanghai AI Lab			2D VAE	T5-XXL	673.68M
Hunyuan-DiT [159]	2024	Tencent Hunyuan			mCLIP & mT5-XL	1.5B	
LuminaT2X [160]	2024	Shanghai AI Lab			LLama2-7B	7B	
Kolors [161]	2024	Kuaishou			ChatGLM3-6B-Base	2.6B	
SD3.0 [80]	2024	Stability AI			CLIP ViT-L & OpenCLIP ViT-bigG & T5-XXL	8B	
Flux.1 [162]	2024	BlackForestLabs			CLIP ViT-L & OpenCLIP ViT-bigG & T5-XXL	12B	
Sora [163]	2024	OpenAI			-	-	
Open-Sora [164]	2024	Hpcaitech			T5-XXL	1.2B	
Open-Sora-Plan [165]	2024	Peking University			3D VAE	T5 & mT5	-
EasyAnimate [166]	2024	Alibaba Group			mCLIP & mT5-XL	1.5B	
CogvideoX [82]	2024	Zhipu AI	T5-XXL	2B/5B			
Moive Gen [83]	2024	Meta	TAE	MetaCLIP & UL2 & ByT5	30B		



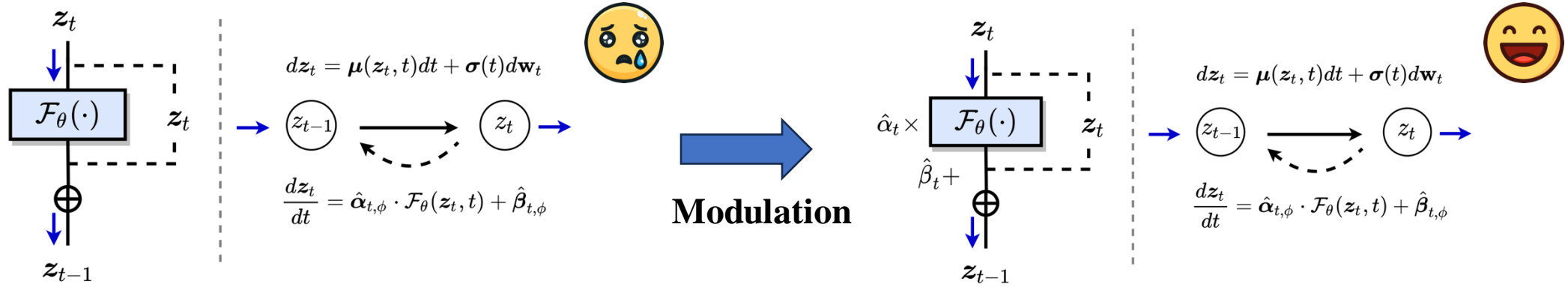
✓ The core of generative intelligence emergence: Scaling Law with increasingly deep stacked networks

# CONTENTS



- 01 Background
- 02 Motivation**
- 03 Methods
- 04 Experiment
- 05 Conclusion



## □ The essential principle of how the generative denoising network works ?



➤ **Core:**

The optimization direction of neural network  $\mathcal{F}_\theta(z_t, t)$   $\equiv$  The direction of inverted diffusion of the data  $z_t$

➤ **Issue:**

1. Asymmetry (coupling) of network predictions:  $\mathcal{F}_\theta(z_t, t)$    $\hat{\alpha}_{t, \phi} \cdot \mathcal{F}_\theta(z_t, t) + \hat{\beta}_{t, \phi}$  **Unbalanced (One Layer)**
2. Training architecture is difficult to scale:  $\mathcal{F}_\theta(z_t, t)$    $\{f_{\theta_1}, \dots, f_{\theta_i}, \dots, f_{\theta_L}\}$  **Unstable (Deep Layer)**

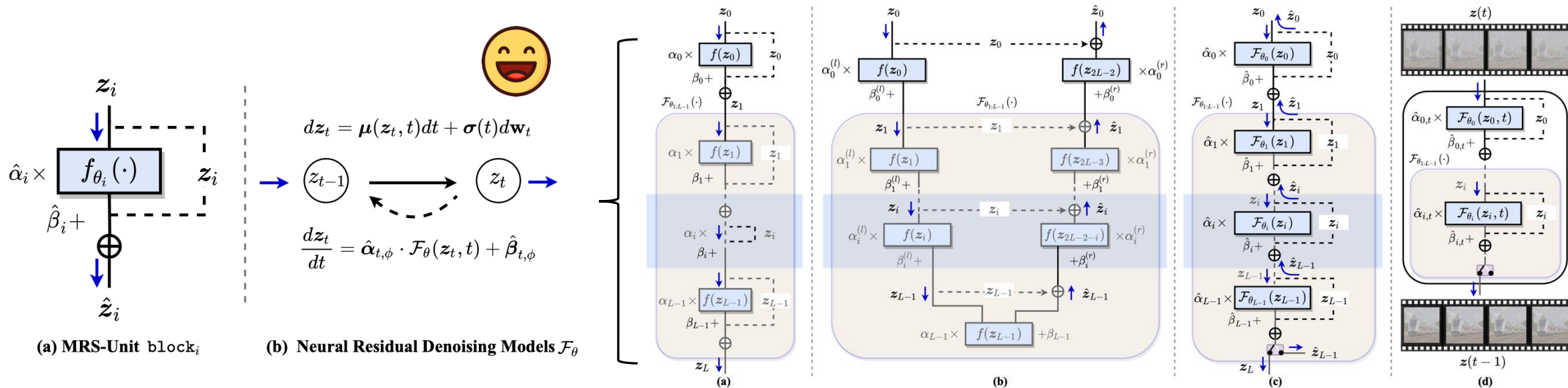
# CONTENTS



- 01 Background
- 02 Motivation
- 03 Methods**
- 04 Experiment
- 05 Conclusion



### Neural Residual Diffusion Models

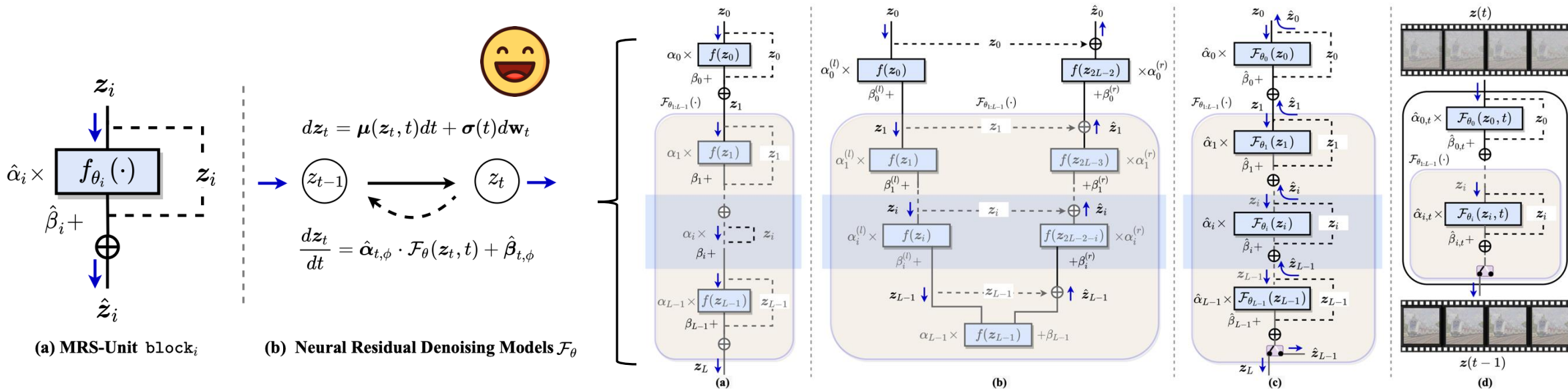


### Gating-Residual Mechanism

$$\begin{cases}
 \hat{z}_i = z_{i+1} = z_i + [\alpha_i \cdot f_{\theta_i}(z_i) + \beta_i] & \text{(Discrete Form)} \\
 \hat{z}_i = \underbrace{\alpha_i^{(l)} \cdot f_{\theta_i^{(l)}}(z_i)}_{\text{read-in branch}} + \underbrace{\alpha_i^{(r)} \cdot f_{\theta_i^{(r)}}(z_{2L-2-i}) + \beta_i^{(r)}}_{\text{read-out branch}} = z_i + \hat{\alpha}_i \cdot \mathcal{F}_{\theta_i}(z_i) + \hat{\beta}_i & \text{(U-shaped)} \quad \text{(Continuous Form)}
 \end{cases}
 \implies \frac{dz_t}{dt} = \hat{\alpha}_\phi \cdot \mathcal{F}_{\theta_t}(z_t) + \hat{\beta}_\phi$$

**Gating Residual Modulation**

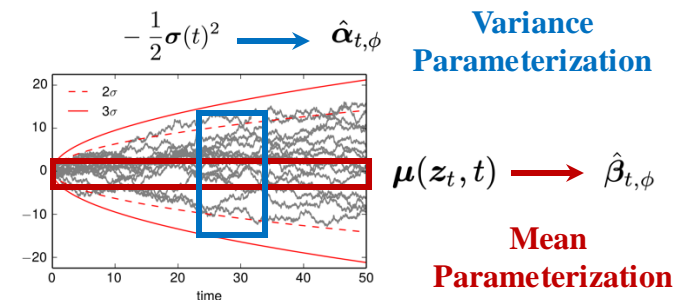
### Neural Residual Diffusion Models



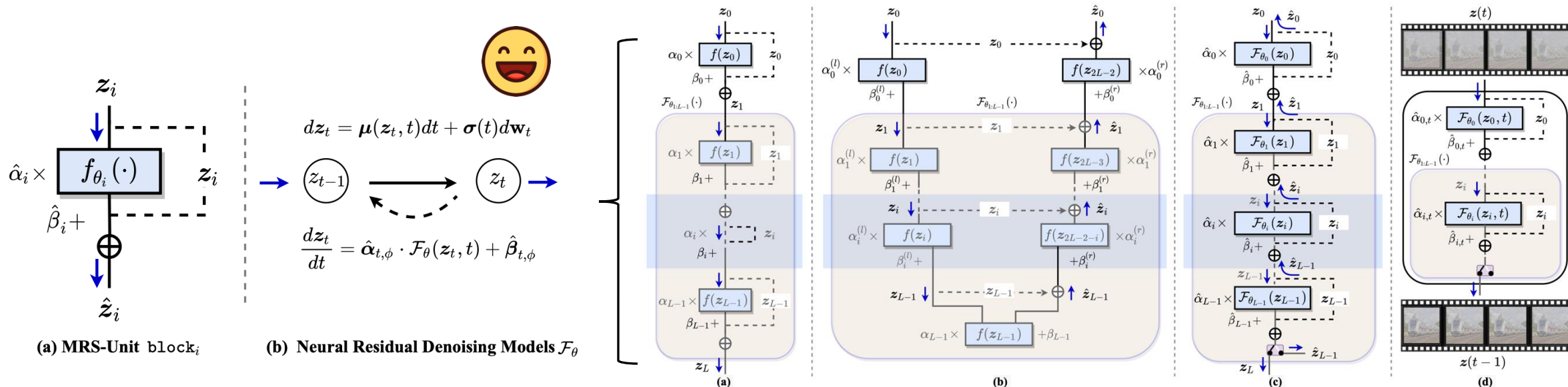
### ➤ Denoising Dynamics Parameterization

$$dz_t = \mu(z_t, t)dt + \sigma(t)dw_t \implies \frac{dz_t}{dt} = \mu(z_t, t) + \sigma(t) \cdot \epsilon_t \quad (\text{Noise-Adding SDE})$$

$$\frac{dz_t}{dt} = \underbrace{\mu(z_t, t)}_{\hat{\beta}_{t,\phi}} - \frac{1}{2}\sigma(t)^2 \cdot \underbrace{\left[ \nabla_z \log p_t(z_t) \right]}_{\hat{\alpha}_{t,\phi}} = \hat{\alpha}_{t,\phi} \mathcal{F}_\theta(z_t, t) + \hat{\beta}_{t,\phi} \quad (\text{Denoising-ODE})$$



## Neural Residual Diffusion Models



### Residual Sensitivity Control

- ◆ To control the **numerical errors** in back-propagation and achieve **steadily and massively scalable training**
- ◆ We introduce **Residual-Sensitivity ODE**:

Define **Residual-Sensitivity** :  $s_t = \frac{d\mathcal{L}}{dz_t} = \frac{d\mathcal{L}}{dz_{t+\delta}} \cdot \frac{dz_{t+\delta}}{dz_t} = s_{t+\delta} \cdot \frac{dz_{t+\delta}}{dz_t}$ .

**Residual-Sensitivity ODE**

$$\frac{ds_t}{dt} = \lim_{\delta \rightarrow 0^+} \frac{s_{t+\delta} - s_t}{\delta} = \lim_{\delta \rightarrow 0^+} \frac{-s_{t+\delta} \cdot \frac{\partial}{\partial z_t} \left( \int_t^{t+\delta} f_\theta(z_t) dt \right)}{\delta} = -s_t \cdot \frac{\partial f_\theta(z_t, t)}{\partial z_t}$$

## ➤ Residual Sensitivity Control

- ◆ To control the **numerical errors** in back-propagation and achieve **steadily and massively scalable training**
- ◆ First, we define **Residual-Sensitivity**:  $\mathbf{s}_t = \frac{d\mathcal{L}}{dz_t}$

$$\mathbf{s}_t = \frac{d\mathcal{L}}{dz_t} = \frac{d\mathcal{L}}{dz_{t+\delta}} \cdot \frac{dz_{t+\delta}}{dz_t} = \mathbf{s}_{t+\delta} \cdot \frac{dz_{t+\delta}}{dz_t} \quad \xrightarrow{\text{(Chain Rule)}} \quad \because dz_{t+\delta} = dz_t + \int_t^{t+\delta} f_\theta(z_t, t) dt. \quad \therefore \mathbf{s}_t = \mathbf{s}_{t+\delta} + \mathbf{s}_{t+\delta} \cdot \frac{\partial}{\partial z_t} \left( \int_t^{t+\delta} f_\theta(z_t, t) dt \right).$$

(Chain Rule) (Euler Solver)

### Residual-Sensitivity ODE

$$\because \mathbf{s}_t = \mathbf{s}_{t+\delta} + \mathbf{s}_{t+\delta} \cdot \frac{\partial}{\partial z_t} \left( \int_t^{t+\delta} f_\theta(z_t, t) dt \right). \quad \therefore \boxed{\frac{ds_t}{dt} = \lim_{\delta \rightarrow 0^+} \frac{\mathbf{s}_{t+\delta} - \mathbf{s}_t}{\delta} = \lim_{\delta \rightarrow 0^+} \frac{-\mathbf{s}_{t+\delta} \cdot \frac{\partial}{\partial z_t} \left( \int_t^{t+\delta} f_\theta(z_t, t) dt \right)}{\delta} = -\mathbf{s}_t \cdot \frac{\partial f_\theta(z_t, t)}{\partial z_t}.}$$

- ◆ Then, we further use the **Euler solver** to obtain the sensitivity  $\mathbf{s}_{t_0}$  :

$$\mathbf{s}_{t_0} = \mathbf{s}_{t_L} + \int_{t_L}^{t_0} \frac{ds_t}{dt} dt = \mathbf{s}_{t_L} - \int_{t_L}^{t_0} \mathbf{s}_t \cdot \frac{\partial f_\theta(z_t, t)}{\partial z_t} dt. \quad \xrightarrow{\text{(non-negativity)}} \quad \mathbf{s}_{t_L} > \mathbf{s}_{t_{L-1}} > \dots > \mathbf{s}_{t_0}$$

(non-negativity) gradually decaying sensitivity !

- ◆ Similarly, we can define **parameter-sensitivity**:  $\mathbf{s}_\theta = \frac{d\mathcal{L}}{d\theta}$ , we can derive:

$$\mathbf{s}_{\theta_0} = \mathbf{s}_{\theta_L} + \int_{t_L}^{t_0} \frac{ds_\theta}{dt} dt = \mathbf{s}_{\theta_L} - \int_{t_L}^{t_0} \mathbf{s}_\theta \cdot \frac{\partial f_\theta(z_t, t)}{\partial \theta} dt. \quad \xrightarrow{\text{(non-negativity)}} \quad \mathbf{s}_{\theta_L} > \mathbf{s}_{\theta_{L-1}} > \dots > \mathbf{s}_{\theta_0}$$



## ➤ Residual Sensitivity Control

◆ So far, we have explored the current situation of the problem:

### *Residual-Sensitivity ODE*

$$\frac{ds_t}{dt} = \lim_{\delta \rightarrow 0^+} \frac{s_{t+\delta} - s_t}{\delta} = \lim_{\delta \rightarrow 0^+} \frac{-s_{t+\delta} \cdot \frac{\partial}{\partial z_t} \left( \int_t^{t+\delta} f_\theta(z_t) dt \right)}{\delta} = -s_t \cdot \frac{\partial f_\theta(z_t, t)}{\partial z_t}.$$

$$\begin{cases} s_{t_0} = s_{t_L} + \int_{t_L}^{t_0} \frac{ds_t}{dt} dt = s_{t_L} - \int_{t_L}^{t_0} s_t \cdot \frac{\partial f_\theta(z_t, t)}{\partial z_t} dt. \\ s_{\theta_0} = s_{\theta_L} + \int_{t_L}^{t_0} \frac{ds_\theta}{dt} dt = s_{\theta_L} - \int_{t_L}^{t_0} s_\theta \cdot \frac{\partial f_\theta(z_t, t)}{\partial \theta} dt. \end{cases}$$

◆ Subsequently, we apply **Gating-Residual** and **Mean-Variance Parameterization** to *Residual-Sensitivity ODE* :

### *Rectified Residual-Sensitivity ODE*

$$\frac{d\hat{s}_t}{dt} = \lim_{\delta \rightarrow 0^+} \frac{\hat{s}_{t+\delta} - \hat{s}_t}{\delta} = -(\alpha_{t,\phi} \cdot \hat{s}_t) \cdot \frac{\partial f_\theta(\hat{z}_t, t)}{\partial \hat{z}_t} - (\beta_{t,\phi} \cdot \hat{s}_t).$$

$$\begin{aligned} \hat{s}_{t_0} &= \hat{s}_{t_L} + \int_{t_L}^{t_0} \frac{d\hat{s}_t}{dt} dt \\ &= \hat{s}_{t_L} - \int_{t_L}^{t_0} \left[ (\alpha_{t,\phi} \cdot \hat{s}_t) \cdot \frac{\partial f_\theta(\hat{z}_t, t)}{\partial \hat{z}_t} + (\beta_{t,\phi} \cdot \hat{s}_t) \right] dt. \end{aligned}$$

◆ Eventually, we can supervise it to achieve *Residual Sensitivity Control* via:

$$\mathcal{L}_s = \|\mathcal{F}_\theta(z_t, t) - \nabla_z \log p_t(z_t)\|_2^2 + \gamma \cdot \sum_L \|\alpha_{t,\phi} \cdot \frac{\partial f_\theta(\hat{z}_t, t)}{\partial \hat{z}_t} - \beta_{t,\phi}\|_2^2 \quad \text{(Rectified Term)}$$



# CONTENTS



01 Background

02 Motivation

03 Methods

**04 Experiment**

05 Conclusion

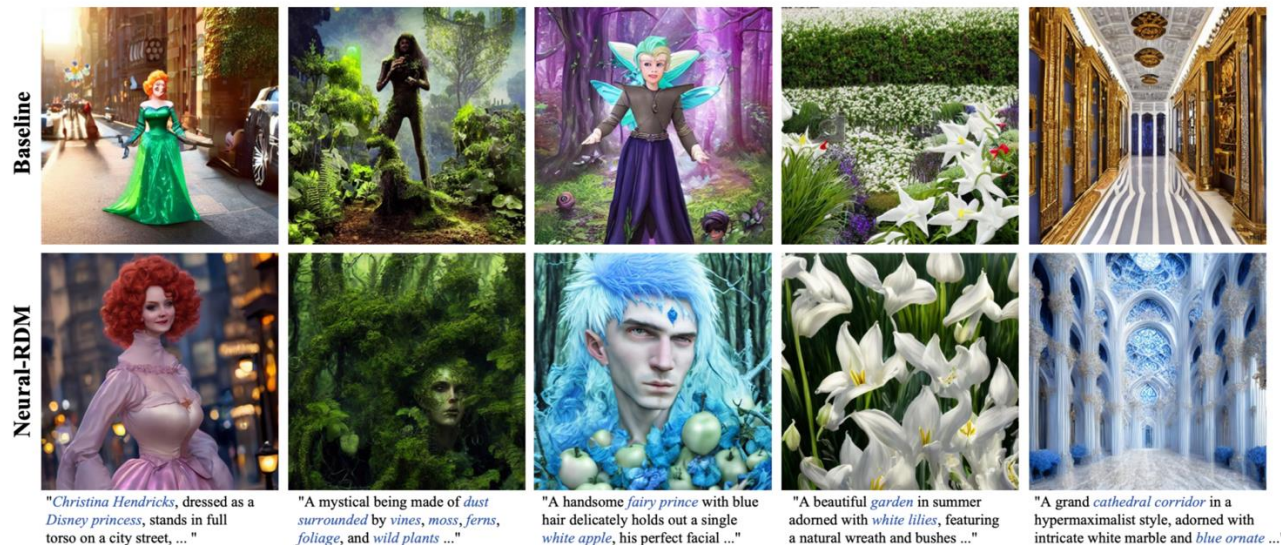


### □ Experiments on Image Synthesis with Deep Scalable Spatial Learning

Architecture	Method	Scalability	Class-to-Image			Text-to-Image		
			FID↓	sFID↓	IS↑	FID↓	sFID↓	IS↑
GAN	BigGAN-deep [56]	✗	6.95	7.36	171.4	-	-	-
	StyleGAN-XL [57]	✗	2.30	4.02	265.12	-	-	-
U-shaped	ADM [58]	✓	10.94	6.02	100.98	-	-	-
	ADM-U	✓	7.49	5.13	127.49	-	-	-
	ADM-G	✓	4.59	5.25	186.70	-	-	-
	LDM-8 [30]	✓	15.51	-	79.03	16.64	11.32	64.50
	LDM-8-G	✓	7.76	-	209.52	9.35	10.02	125.73
	LDM-4	✓	10.56	-	103.49	12.37	11.58	94.65
F-shaped	LDM-4-G	✓	3.60	-	247.67	3.78	5.89	182.53
	DiT-XL/2 [59]	✓	9.62	6.85	121.50	8.53	5.47	144.26
	DiT-XL/2-G	✓	2.27	4.60	278.24	3.53	5.48	175.63
Unified	Latte-XL [60]	✓	2.35	5.17	224.75	2.74	5.35	195.03
	Neural-RDM-U (Ours)	✓✓	3.47	5.08	256.55	2.25	4.36	235.35
	Neural-RDM-F (Ours)	✓✓	2.12	3.75	295.32	2.46	5.65	206.32

Table 1: The main results for image generation on ImageNet [61] (Class-to-Image) and JourneyDB [53] (Text-to-Image) with  $256 \times 256$  image resolution. We highlight the best value in blue, and the second-best value in green. The Scalability column indicates the scaling capability of the parameter scale and architecture.

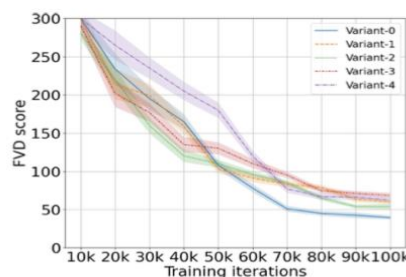
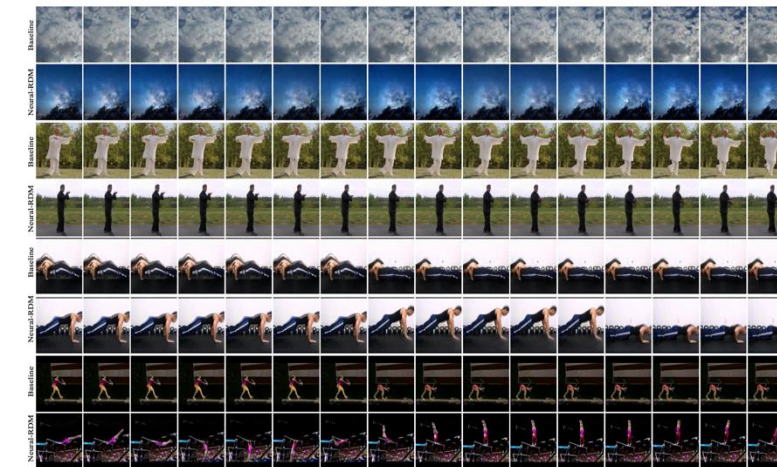
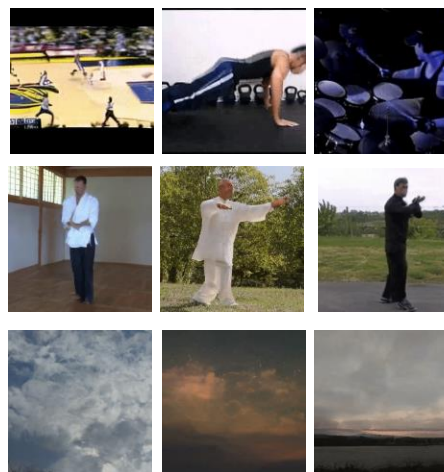
- ◆ Neural-RDMs have obtained competitive and state-of-the-art results across image synthesis benchmarking.
- ◆ Benefiting from the rectification of generative dynamics, it highlights the semantics of the subject more.



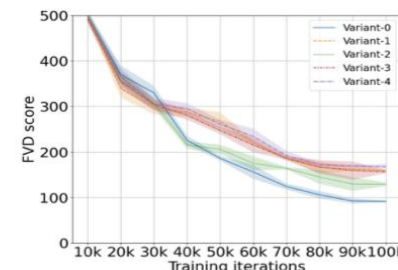
## □ Experiments on Video Generation with Deep Scalable Temporal Learning

Method	Scalability	Frame Evaluation		None-to-Video		Class-to-Video
		FID↓	IS↑	SkyTimelapse (FVD↓)	Taichi-HD (FVD↓)	UCF-101 (FVD↓)
MoCoGAN [71]	✗	23.97	10.09	206.6	-	2886.9
MoCoGAN-HD [72]	✗	7.12	23.39	164.1	128.1	1729.6
DIGAN [73]	✗	19.10	23.16	83.11	156.7	1630.2
StyleGAN-V [70]	✗	9.45	23.94	79.52	-	1431.0
MoStGAN-V [74]	✗	-	-	65.30	-	1380.3
PVDM [75]	✓	29.76	60.55	75.48	540.2	1141.9
LVDM [12]	✓	-	-	95.20	99.0	<b>372.0</b>
VideoGPT [76]	✓	22.70	12.61	222.7	-	2880.6
Latte-XL [60]	✓	5.02	68.53	59.82	159.60	477.97
<b>Neural-RDM (Ours)</b>	✓✓	<b>3.35</b>	<b>85.97</b>	<b>39.89</b>	<b>91.22</b>	<b>461.03</b>

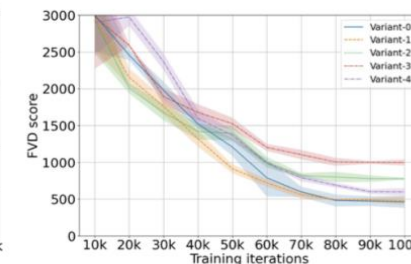
Table 2: The main results for video generation on the SkyTimelapse [62], Taichi-HD [63] and UCF-101 [64] with  $256 \times 256$  resolution of each frame. We highlight the best value in blue, and the second-best value in green.



(a)



(b)



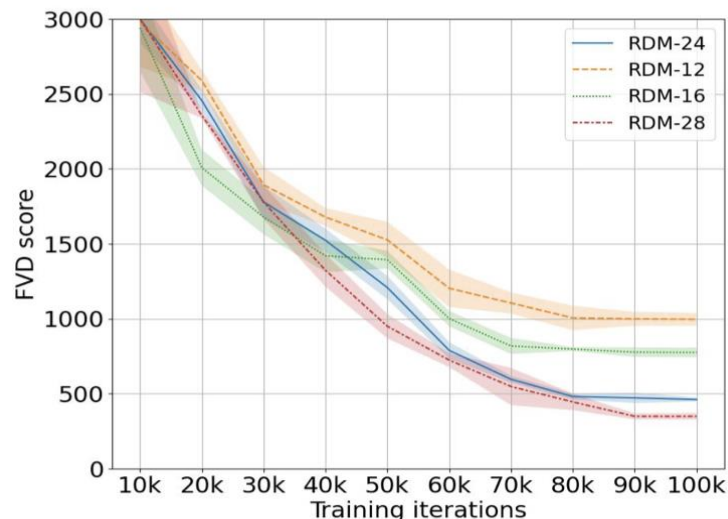
(c)

Figure 6: (a), (b), and (c) respectively illustrate the performance of the five residual structures variant models across the SkyTimelapse [62], Taichi-HD[63], and UCF-101 [64].

- ◆ Neural-RDMs (flow-shaped version) basically achieves the best results, except for the second-best results in class-to-video evaluation.
- ◆ Compare with the baseline, Neural-RDM maintains temporal coherence and consistency, resulting in smoother and more dynamic video frames.



## Comparison Experiments of Gating Residual Variants and Deep Scalability



- ① *Variant-0 (Ours)*:  $z_{i+1} = z_i + \alpha f(z_i) + \beta$ ;
- ② *Variant-1 (AdaLN [77])*:  $z_{i+1} = z_i + f(\alpha z_i + \beta)$ ;
- ③ *Variant-2*:  $z_{i+1} = \alpha z_i + f(z_i) + \beta$ ;
- ④ *Variant-3 (ResNet [78])*:  $z_{i+1} = z_i + f(z_i)$ ;
- ⑤ *Variant-4 (ReZero [79])*:  $z_{i+1} = z_i + \alpha f(z_i)$ .

- ◆ As the number of training steps increases, almost all variants can converge effectively, but only Variant-0 (Our approach) achieves the best FVD scores.
- ◆ As the depth of residual units increases, the performance of the model can be further improved, which further highlights the deep scalability advantage of Neural-RDM.

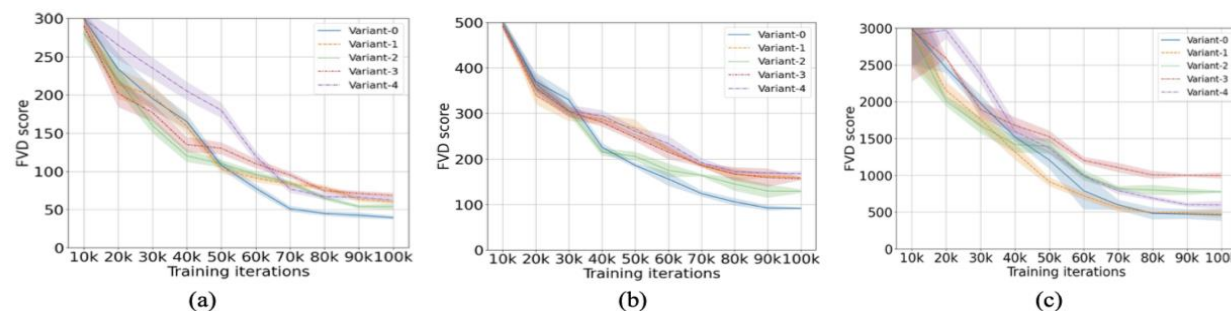


Figure 6: (a), (b), and (c) respectively illustrate the performance of the five residual structures variant models across the SkyTimelapse [62], Taichi-HD[63], and UCF-101 [64].

# CONTENTS



01 Background

02 Motivation

03 Methods

04 Experiment

**05 Conclusion**

- ✓ **Propose a unified neural residual diffusion models framework**

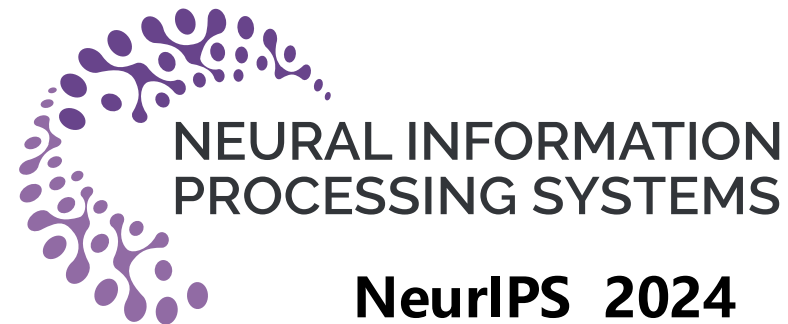
We practically unify u-shaped and flow-shaped stacking networks and to propose a unified and deep scalable neural residual diffusion model framework.

- ✓ **Parameterize the mean-variance scheduler for excellent dynamics consistency**

Moreover, we theoretically parameterize the previous human-designed mean-variance scheduler and demonstrate excellent dynamics consistency.

- ✓ **Adequate and extensive experiments and analyses**

Experimental results on various generative tasks show that Neural-RDM obtains the best results, and extensive experiments also demonstrate the advantages in improving the fidelity, consistency of generated content and supporting large-scale scalable training.



# Thanks for your listening

**Zhiyuan Ma**

**(<https://ponymzy.github.io>)**

Department of Electronic Engineering,

Tsinghua University, Beijing, China

[mzyth@tsinghua.edu.cn](mailto:mzyth@tsinghua.edu.cn)