

# Jump Your Steps: Optimizing Sampling Schedule of Discrete Diffusion Models

---

Yonghyun Park<sup>1</sup>, Chieh-Hsin Lai<sup>1</sup>, Satoshi Hayakawa<sup>2</sup>, Yuhta Takida<sup>1</sup>, Yuki Mitsufuji<sup>1,2</sup>

<sup>1</sup>SONY AI, <sup>2</sup>Sony Group Corporation

# Discrete Diffusion Models (DDMs)

- Recently, DDMs show promising results on in various fields:
  - Image
  - Language
  - Gene, Protein

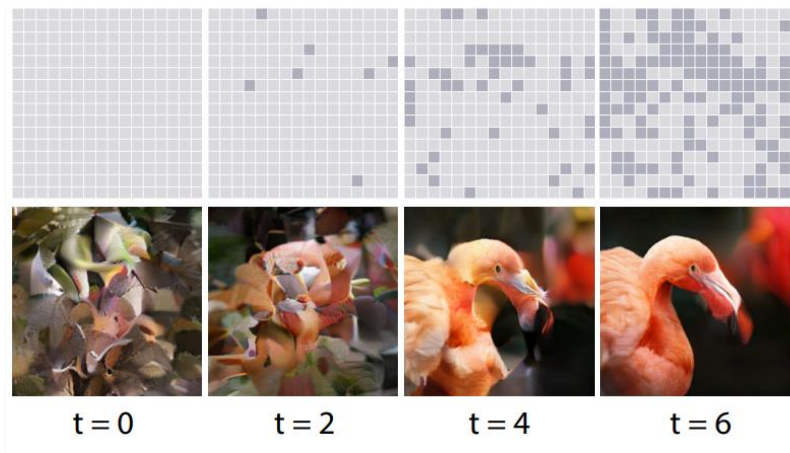
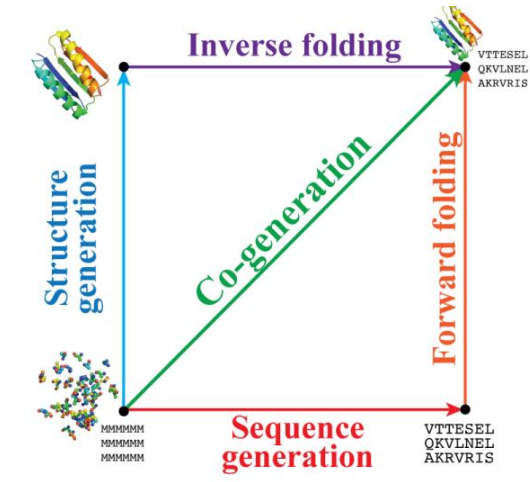


Image Generation  
(MaskGIT)



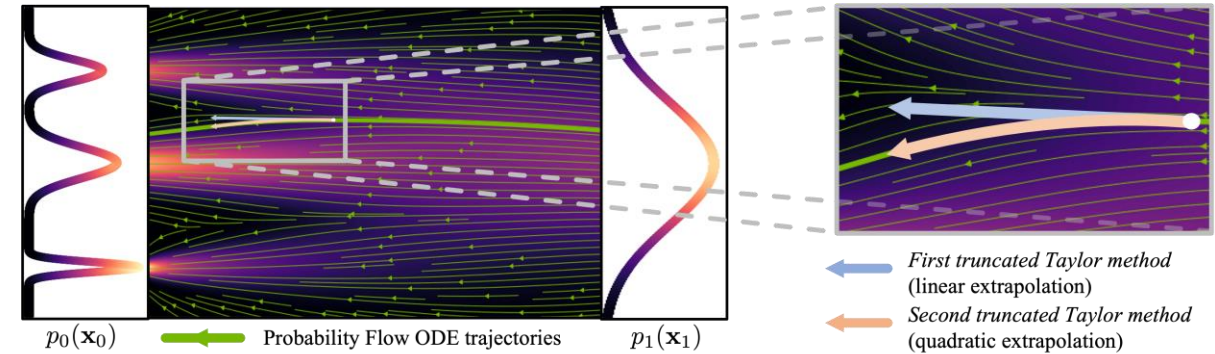
Language Generation  
(LLaDA)



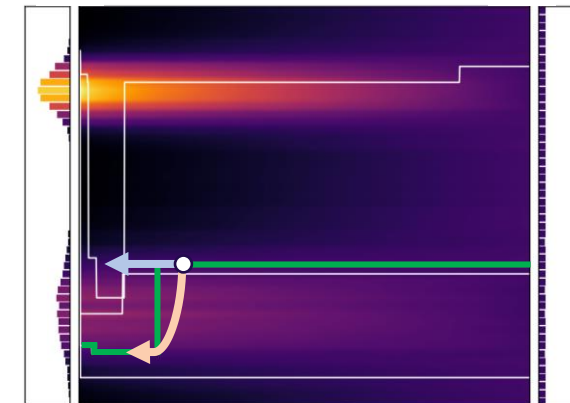
Gene / Protein Generation  
(DFM)

# Challenges with DDMs

- **Slow sampling:**
  - DMs require multiple NFEs for generation.
- **Continuous DMs**
  - Solution: high-order ODE / SDE solver.
  - There are many previous work for fast and reliable differential equation solver.
- **Discrete DMs**
  - Solution: high-order Jump process solver?
  - There are little work for efficient Markov jump process solver.



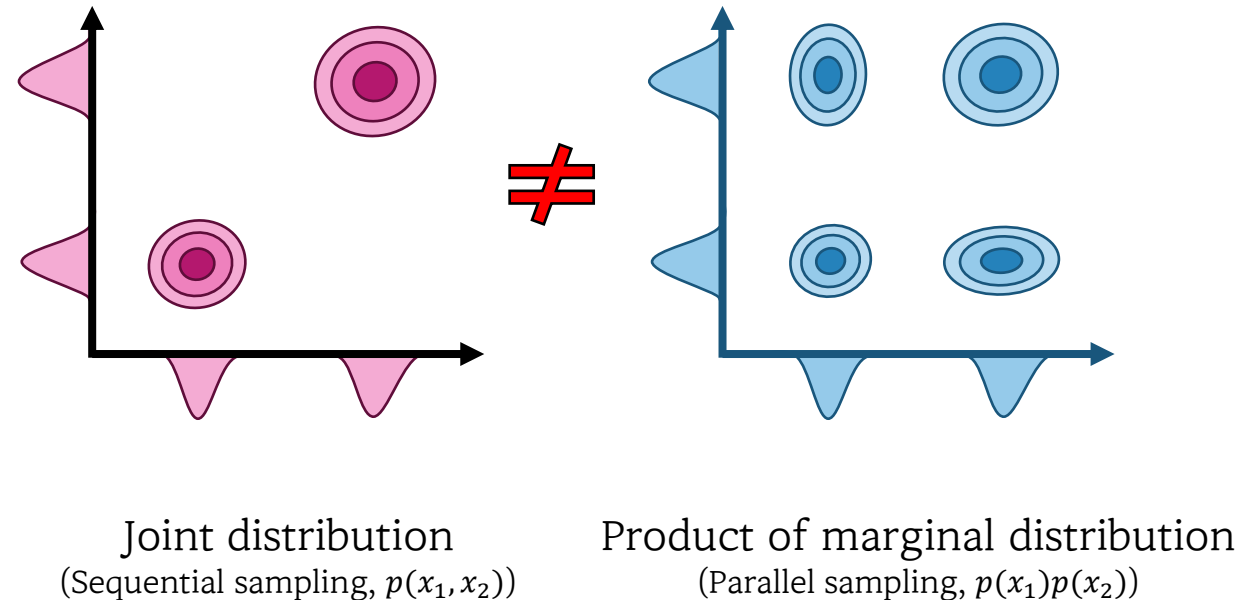
**GENIE: Higher-Order Denoising Diffusion Solvers**



**High-order Markov Jump process solvers are difficult to develop.**

# Challenges with DDMs

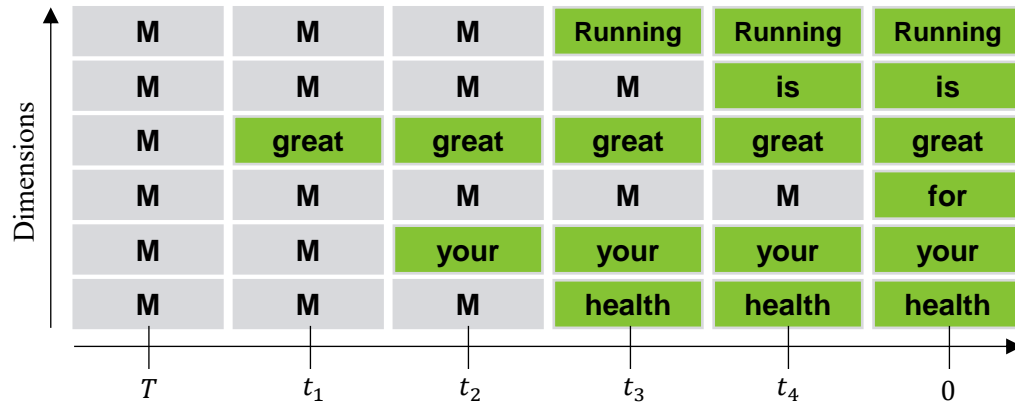
- **Faster sampling for DDMs: Parallel sampling**
  - Sampling multiple tokens within single step
    - $k$ -Gillespie: Sample  $k$  tokens simultaneously.
    - $\tau$ -leaping: Sample all tokens simultaneously within a given time frame.
- **Compounding Decoding Error (CDE)**
  - Generally, parallel sampling leads to error.
    - Parallel sampling:  $p(x_1)p(x_2)$
    - Sequential sampling:  $p(x_1)p(x_2|x_1)$
  - $\text{CDE} := D_{KL}(p(x_1, x_2) || p(x_1)p(x_2))$



# Main Idea: Motivating Example

- CDE depends on timesteps

*Discrete sequence trajectory:  $(X_t)_{t \in [0, T]}$*

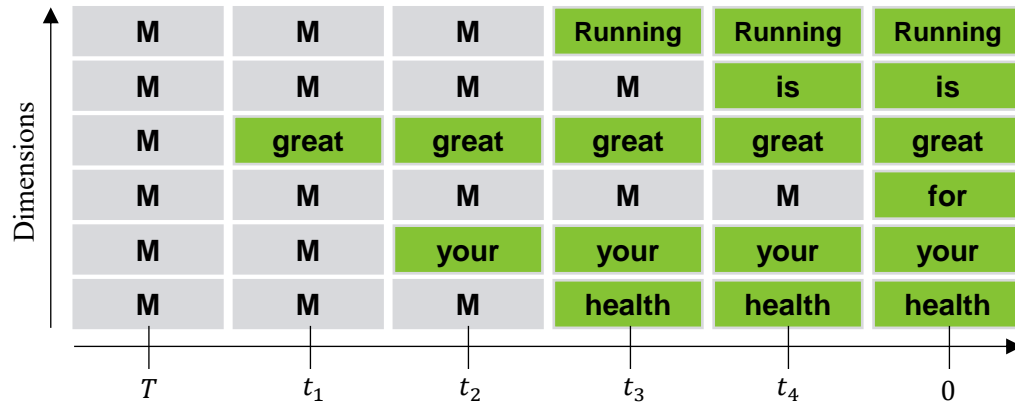


Ground-truth sampling trajectory

# Main Idea: Motivating Example

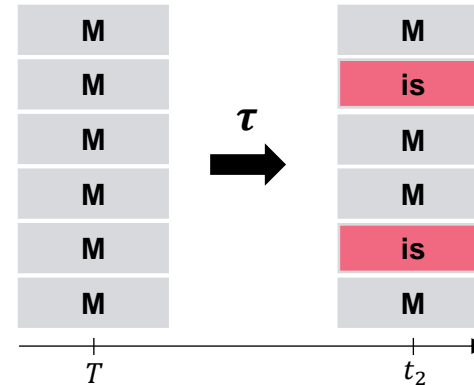
- CDE depends on timesteps

Discrete sequence trajectory:  $(X_t)_{t \in [0, T]}$



Ground-truth sampling trajectory

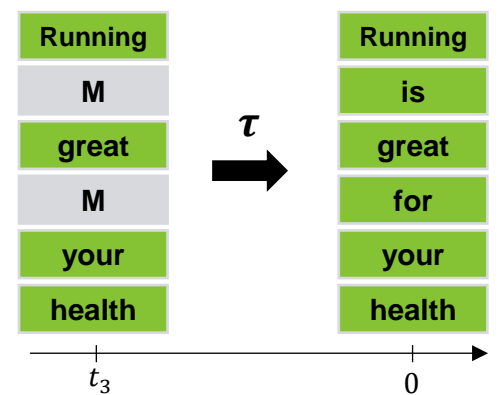
$t \approx T$



**High CDE**

Parallel sampling  
hurts performance.

$t \ll T$



**Low CDE**

Minimal impact from  
parallel generation.

# Main Idea: Jump Your steps

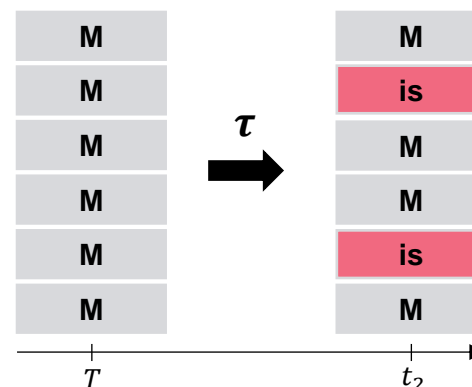
- CDE depends on timesteps

Discrete sequence trajectory:  $(X_t)_{t \in [0, T]}$



Ground-truth sampling trajectory

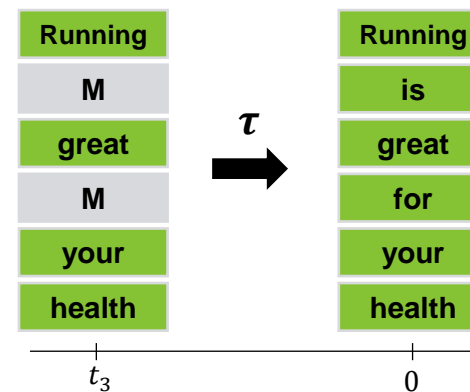
$t \approx T$



**High CDE**

Parallel sampling hurts performance.

$t \ll T$



**Low CDE**

Minimal impact from parallel generation.

Optimizing Sampling Schedule of DDMs

Use small step size  
when decoding error is large

Use large step size  
when decoding error is small

# Theory

Theorem 3.1

Eq. (10)

$$\mathcal{D}_{\text{KL}}(\mathbb{P}_0 \| \mathbb{Q}_0^{t_0 \rightarrow t_1 \rightarrow \dots \rightarrow 0}) \leq \sum_{i=0}^{N-1} \mathcal{D}_{\text{KL}}(\mathbb{P}_{t_{i+1}} \| \mathbb{Q}_{t_{i+1}}^{t_i \rightarrow t_{i+1}}) \leq \mathcal{D}_{\text{KL}}(\mathbb{P}_{\text{paths}} \| \mathbb{Q}_{\text{paths}}^{t_0 \rightarrow t_1 \rightarrow \dots \rightarrow 0})$$

Sampling Quality

Eqs. (3, 5)

Theorem 3.2

$$\sum_{i=0}^{N-1} \mathcal{E}_{\text{CDE}}(t_i \rightarrow t_{i+1})$$

$$\text{KLUB}(\mathbb{P}_0 \| \mathbb{Q}_0^{t_0 \rightarrow t_1 \rightarrow \dots \rightarrow 0})$$

Motivation: Minimize CDE

Q: Why should we minimize CDE?

A: It improves the sampling quality.



# Theory

Theorem 3.1

Eq. (10)

$$\mathcal{D}_{\text{KL}}(\mathbb{P}_0 \| \mathbb{Q}_0^{t_0 \rightarrow t_1 \rightarrow \dots \rightarrow 0}) \leq \sum_{i=0}^{N-1} \mathcal{D}_{\text{KL}}(\mathbb{P}_{t_{i+1}} \| \mathbb{Q}_{t_{i+1}}^{t_i \rightarrow t_{i+1}}) \leq \mathcal{D}_{\text{KL}}(\mathbb{P}_{\text{paths}} \| \mathbb{Q}_{\text{paths}}^{t_0 \rightarrow t_1 \rightarrow \dots \rightarrow 0})$$

|  
| Eqs. (3, 5)  
|

$$\sum_{i=0}^{N-1} \mathcal{E}_{\text{CDE}}(t_i \rightarrow t_{i+1})$$

**Motivation: Minimize CDE**

|  
| Theorem 3.2  
|

$$\text{KLUB}(\mathbb{P}_0 \| \mathbb{Q}_0^{t_0 \rightarrow t_1 \rightarrow \dots \rightarrow 0})$$

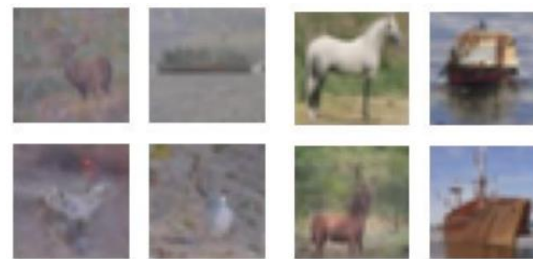
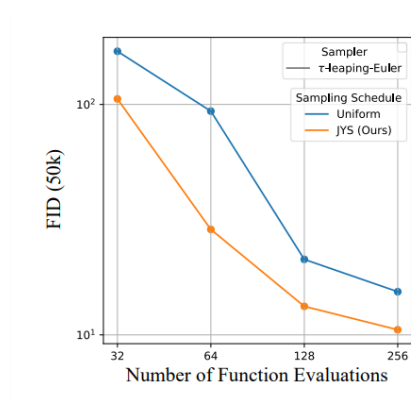
**Trackable Upper-bound**

Q: How could we minimize CDE?

A: Using trackable upper-bound.

# Results

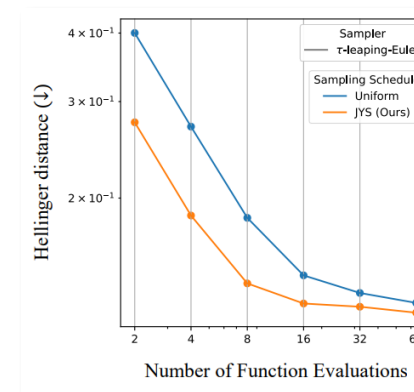
- Improved Sampling quality:
  - JYS sampling scheduler achieves performance improvement regardless of data domain or model.



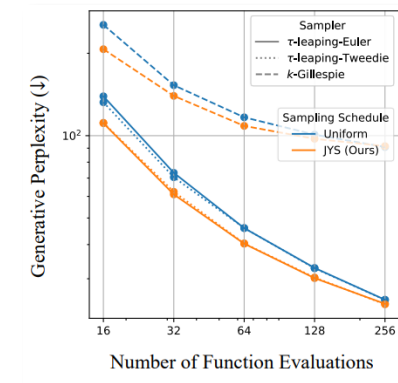
Uniform

Ours

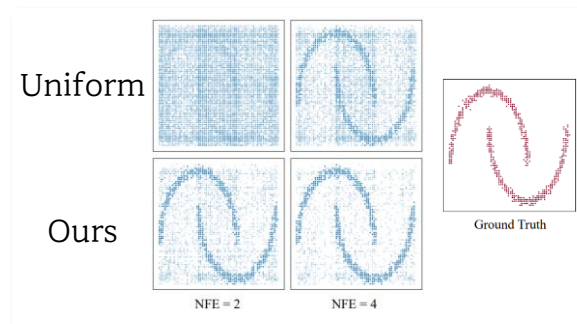
CIFAR10



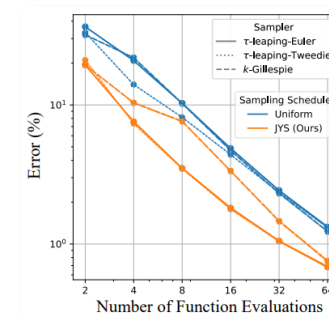
Piano



Language (GPT-2 scale)



2D Dataset



CountDown

Thank you