



SCHOOL OF **COMPUTING** & DATA SCIENCE The University of Hong Kong

# On the Sampling Theory for Auto-Regressive **Diffusion Inference**

May. 28-29, 2025 HKU IDS Interdisciplinary Workshop – Exploring the Foundations: Fundamental AI and Theoretical Machine Learning Joint work with Xunpeng Huang, Yujin Han, Hanze Dong, Yi Zhang, Yian Ma, and Tong Zhang



### **HKU Musketeers Foundation Institute of Data Science** 香港大學同心基金數據科學研究院

# **Difan Zou**



### **Diffusion Model**

Sohl-Dickstein et al, ICML 2015

### Adding noise via OU process

 $p_0(x)$ 



 $p_T^{\leftarrow}(x) \approx p_0(x)$ 



Denoising via reversing OU process





 $p_0^{\leftarrow}(x) = N(0,\mathbf{I})$ 



- the discrete space.
- next-token prediction using transformer models.

https://deepgenerativemodels.github.io/notes/autoregressive/

### The data is decomposed into $x = [x_1, x_2, \dots, x_n]$ , each token is encoded in

> The tokens are generated in a sequential manner:  $x_k \sim p(x | x_1, \dots, x_{k-1})$ In the standard AR generative model,  $p(\cdot | x_1, \dots, x_{k-1})$  is implemented as



### **AR Diffusion Inference**



- Vhen modeling the conditional distribution  $p(\cdot | x_1, \dots, x_{k-1})$ , using **discrete** encoding may lead to information loss.
- AR diffusion inference leverages diffusion model to model  $p(\cdot | x_1, \dots, x_{k-1})$ , which enables conditional generation in continuous space.

Li et al., Autoregressive Image Generation without Vector Quantization, NeurIPS 2024

### **AR Diffusion Inference**

- Procedure of AR diffusion model:
  - Decompose the data x into K parts.
  - Given the previous tokens  $x_1, \ldots, x_k$  (in continuous space), we first leverage a encoder  $g_{\theta_a}$ to get the condition embedding z =
  - $x_{k+1} = s_{\theta}(\xi, k, z).$

$$g_{\theta_e}(x_{[1:k]}).$$

Then, we leverage the conditional diffusion model  $s_{\theta}$  to generate the next token



# **AR Diffusion Model VS. Standard Diffusion Model**

- Some facts:
  - > When K = 1, AR DM reduces to the standard DM.
  - AR DM learn many conditional probabilities, while standard DM directly learns the joint distribution.
  - AR DM needs to perform inference in a sequential manner (token-by-token), while standard DM generates data in parallel.
  - Each step of AR DM operates on the space with low dimension, while DM works on full dimension space.

Goal of our work: build the theoretical foundation of AR DM to explain the pros and cons of the auto-regressive paradigm.



# **Sampling Error of AR Diffusion Model**

Inference Process



Sampling error analysis:

Our analysis needs to (1) study single-step sampling error; and (2) study the error propagation during AR process.

In each step, the DM sampling process  $\hat{p}_{\theta}(\cdot | x_1, \dots, x_k)$  could lead to sampling error. The single-step error will accumulate during the sequential sampling process.

# **DM Sampling Algorithm: DDPM**

Forward process: adding noise

 $p_{t|0}(\mathbf{x}_t \,|\, \mathbf{x}_0)$ 

Reverse process: generating via denoising  $p_{t|t-1}(\mathbf{x}_{t-1} | \mathbf{x}_t) \approx N(\mu(\mathbf{x}_t))$ 



Data distribution



= 
$$N(\alpha_t \mathbf{x}_0, \beta_t \mathbf{I})$$
 Gaussian Transition Kernel

$$t, t), \Sigma(\mathbf{x}_t, t))$$

Use a gaussian kernel to approximate the single reverse step.





### **Theoretical Analysis for DDPM**

### Error Decomposition of DDPM

$$d\mathbf{x}_{t}^{\leftarrow} = \left(\mathbf{x}_{t}^{\leftarrow} + \mathbf{v}_{\theta}(\mathbf{x}_{k\eta}^{\leftarrow}, k)\right) dt + \sqrt{2} dB_{t}, \quad t \in [k\eta, (k+1)\eta]$$
• Truncation error:  $p(\mathbf{x}_{T}) \approx N(0, \mathbf{I})$ 
• Discretization error:  $\nabla \ln p(\mathbf{x}_{k\eta}) \approx \nabla p_{t}(\mathbf{x})$ 
• Score estimation error:  $\mathbf{v}_{\theta}(\mathbf{x}) \approx \nabla p_{T-t}(\mathbf{x})$ 
• Error propagation

 $\mathrm{d}\hat{\mathbf{x}}_t = (\hat{\mathbf{x}}_t + 2\nabla \ln p_{T-t}(\hat{\mathbf{x}}_t)) \mathbf{c}$ 



$$\mathrm{d}t + \sqrt{2}\mathrm{d}B_t, \quad \hat{\mathbf{x}}_0 \sim p_T.$$

### **Theoretical Analysis for DDPM**

To guarantee small sampling error

Discretization error:  $\|\nabla \ln p(\mathbf{x}_{k\eta}) - \nabla p_t(\mathbf{x})\| \sim \operatorname{poly}(\eta)$ 



Chen, S., Chewi, S., Li, J., Li, Y., Salim, A., & Zhang, A. Sampling is as easy as learning the score: theory for diffusion models with minimal data assumptions. In The Eleventh International Conference on Learning Representations. Lee, H., Lu, J., & Tan, Y. (2022). Convergence for score-based generative modeling with polynomial complexity. Advances in Neural Information Processing *Systems*, *35*, 22870-22882.



### • Consequently, we need to use a large iteration number $K = T/\eta \sim \log(1/\epsilon)/\eta$ .

### **Revisit DDPM Algorithm**



DDPM algorithm aims to sample  $\mathbf{x}_{t-1}$  given  $\mathbf{x}_t$  using a Gaussian kernel.  $p_{t-1|t}(\mathbf{x}_{t-1} \mid \mathbf{x}_t) \approx$ 

- Gaussian kernel is easy to implement.
- Gaussian approximation can only be good when  $\eta$  is small.
- Then we must require many reverse sampling steps  $(T/\eta)$ .
- What if we consider larger  $\eta$ ?

$$\approx N(\mu(\mathbf{x}_t, t), \boldsymbol{\Sigma}(\mathbf{x}_t, t))$$



We can get the exact formula of the **reverse transition kernel**  $p_{t|t'}(\mathbf{x}_t | \mathbf{x}_{t'})$ 

$$p_{t|t'}(\mathbf{x}_t \,|\, \mathbf{x}_{t'}) \propto p_{t|t'}(\mathbf{x}_t \,|\, \mathbf{x}_{t'}) \cdot p_{t'}(\mathbf{x}_{t'}) = p_{t'|t}(\mathbf{x}_{t'} \,|\, \mathbf{x}_{t}) \cdot p_{t}(\mathbf{x}_{t'})$$

Forward transition kernel, which is a Gaussian



Marginal distribution at time *t*,  $\propto e^{-f_t(\mathbf{x}_t)}$ 

$$-f_t(\mathbf{x}) - \frac{\|\mathbf{x}' - \mathbf{x} \cdot e^{-2(t'-t)}\|^2}{2(1 - e^{-2(t'-t)})}\right)$$



- When t' t is very small, the quadratic term dominates, the kernel can be well approximated by Gaussian, i.e., DDPM.
- When t' t is large, we need to use more complicated sampler to achieve high-accuracy sampling.

# **Trade-off of RTK** $\mathbf{X}_{t-1} \quad \mathbf{X}_t$

### **Consider fixed** segment

- subproblems increases (e.g., DDPM).
- becomes easier.
- How to pick a proper  $\eta$ ?



• Smaller  $\eta$  implies easier sampling subproblems, but the number of

Larger  $\eta$  implies smaller number of subproblems, but the problem

### Hardness of RTK

 $p_{(k-1)\eta|k\eta}(\mathbf{x} \,|\, \mathbf{x}') = \exp$ **RTK** target

We can verify the **log-concavity**, let's check the Hessian of  $\log p_{(k-1)\eta|k\eta}(\mathbf{x} \mid \mathbf{x}')$ 

$$-\nabla_{\mathbf{x}}^{2}\log p_{(k-1)\eta|k\eta}(\mathbf{x} \,|\, \mathbf{x}') = -\nabla_{\mathbf{x}}^{2}f_{(k-1)\eta}(\mathbf{x}) + \frac{e^{-2\eta}}{1 - e^{2\eta}}$$

- Assume that  $f_t(\mathbf{x})$  is L-smooth.
- Then we can set  $\eta \leq 1/2 \cdot \log(1 + 1/(2L))$  to ensure the log-concavity.

$$p\left(-f_{(k-1)\eta}(\mathbf{x}) - \frac{\|\mathbf{x}' - \mathbf{x} \cdot e^{-2\eta}\|^2}{2(1 - e^{-2\eta})}\right)$$

• We can just set  $\eta = 1/2 \cdot \log(1 + 1/(2L))$  to ensure a small number of subproblems.

# **Diffusion Inference with RTK**

### **Algorithm** INFERENCE WITH REVERSE TRANSITION KERNEL (RTK)

- size  $\eta$ , required convergence accuracy  $\epsilon$ ;
- 2: for k = 0 to K 1 do
- 3: transition kernel, i.e.,

$$\boldsymbol{\rho}_{(k+1)\eta|k\eta}^{\leftarrow}(\boldsymbol{z}|\hat{\boldsymbol{x}}_{k\eta}) \propto \exp\left(-g(\boldsymbol{z})\right) \coloneqq \exp\left(-f_{(K-k-1)\eta}(\boldsymbol{z}) - \frac{\|\hat{\boldsymbol{x}}_{k\eta} - \boldsymbol{z} \cdot \boldsymbol{e}^{-\eta}\|^2}{2(1 - \boldsymbol{e}^{-2\eta})}\right).$$
(23)

4: **end for** 

### **Efficiently solving RTK is crucial.**

5: return  $\hat{\mathbf{x}}_{K}$ .



1: Input: Initial particle  $\hat{\mathbf{x}}_0$  sampled from the standard Gaussian distribution, Iteration number K, Step

Draw sample  $\hat{\mathbf{x}}_{(k+1)\eta}$  with MCMCs from  $\hat{p}_{(k+1)\eta|k\eta}(\cdot|\hat{\mathbf{x}}_{k\eta})$  which is closed to the ground-truth reverse

# Solving RTK

. . .

Sampling from the distribution  $p_{(k-1)\eta|k\eta}$ 

Multiple sampling algorithms can be applied:

- Unadjusted Langevin Algorithms (ULA)
- Underdamped Langevin Dynamics (ULD)
  - provide fast sampling and requires mild assumptions
- Metropolis Adjusted Langevin Algorithms (MALA)
  - Provide even faster sampling but need stricter assumptions
- Hamiltonian Monte Carlo (HMC)

$$(\mathbf{x} | \mathbf{x}') = \exp\left(-f_{(k-1)\eta}(\mathbf{x}) - \frac{\|\mathbf{x}' - \mathbf{x} \cdot e^{-2\eta}\|^2}{2(1 - e^{-2\eta})}\right)$$

# **RTK** with ULD

### Algorithm 3 ULD FOR RTK INFERENCE

- $\mathcal{N}(\mathbf{0}, e^{2\eta} 1) \otimes \mathcal{N}(\mathbf{0}, \mathbf{I});$
- 3: for t = s to S 1 do
- Draw noise sample pair  $(\xi_s^z, \xi_s^v)$  from a Gaussian type distribution. 4:
- $\hat{\boldsymbol{z}}_{s+1} = \hat{\boldsymbol{z}}_s + \gamma^{-1} (1 e^{-\gamma\tau}) \hat{\boldsymbol{v}}_s \gamma^{-1} (\tau \gamma^{-1} (1 e^{-\gamma\tau})) s_{\theta}(\hat{\boldsymbol{z}}_s) + \xi_s^z$ 5:
- $\hat{\boldsymbol{v}}_{s+1} = e^{-\gamma\tau} \hat{\boldsymbol{v}}_s \gamma^{-1} (1 e^{-\gamma\tau}) s_{\theta}(\hat{\boldsymbol{z}}_s) + \xi_t^v$ 6:
- 7: return  $\boldsymbol{z}_S$ ;

Here the random Gaussian vectors are generated via:

$$(\xi_s^z, \xi_s^v) \sim \mathcal{N} \left( \mathbf{0}, \begin{bmatrix} \frac{2}{\gamma} \left( \tau - \frac{2}{\gamma} \left( \mathbf{1} - e^{-\gamma \tau} \right) \right) + \frac{1}{2\gamma} \left( \mathbf{1} - e^{-2\gamma \tau} \right) & \frac{1}{\gamma} \left( \mathbf{1} - 2e^{-\gamma \tau} + e^{-2\gamma \tau} \right) \\ \frac{1}{\gamma} \left( \mathbf{1} - 2e^{-\gamma \tau} + e^{-2\gamma \tau} \right) & \mathbf{1} - e^{-2\gamma \tau} \end{bmatrix} \right)$$

1: Input: Returned particle of the previous iteration  $x_0$ , current iteration number k, inner iteration number S, inner step size  $\tau$ , velocity diffusion coefficient  $\gamma$ , required convergence accuracy  $\epsilon$ ; 2: Initialize the particle and velocity pair, i.e.,  $(\hat{z}_0, \hat{v}_0)$  with a Gaussian type product measure, i.e.,

> Neural network approximation of the score function  $\nabla f_{\rm s}(\hat{\mathbf{z}}_{\rm s})$

### **RTK with ULD: Rates**

**Theorem:** Under the assumptions on the score estimation

$$\eta = 1/2 \cdot \log[(2L+1)/2L]$$
 and  $K = 4L \cdot \log[((1+L^2)d + ||\nabla f_*(\mathbf{0})||^2)^2 \cdot \epsilon^{-2}]$ 

and implement Step 3 of Alg. 2 with ULD. For the k-th run of ULD, we require Gaussian-type initialization and high-accurate score estimation, i.e.,

$$\hat{\pi}_0 = \mathcal{N}(\mathbf{0}, e^{2\eta} - 1) \otimes \mathcal{N}(\mathbf{0}, \mathbf{I}) \text{ and } \epsilon_{\text{score}} = \tilde{\mathcal{O}}(\epsilon/\sqrt{L}).$$

If we set the hyperparameters of inner loops as follows. the step size and the iteration number as

$$\tau = \tilde{\Theta} \left( \epsilon d^{-1/2} L^{-1/2} \cdot \left( \log \left[ \frac{L(d + m_2^2 + \|\boldsymbol{x}_0\|^2)}{\epsilon^2} \right] \right)^{-1/2} \right)$$
$$S = \tilde{\Theta} \left( \epsilon^{-1} d^{1/2} \cdot \left( \log \left[ \frac{L(d + m_2^2 + \|\boldsymbol{x}_0\|^2)}{\epsilon^2} \right] \right)^{1/2} \right).$$

It can achieve  $\|\hat{p}_{K\eta} - p_*\|_{TV} \lesssim \epsilon$  with an  $\tilde{\mathcal{O}}(L^2 d^{1/2} \epsilon^{-1})$  gradient complexity.

### **RTK with MALA**

### Algorithm 2 MALA/PROJECTED MALA FOR RTK INFERENCE

- number S, inner step size  $\tau$ , required convergence accuracy  $\epsilon$ ;
- 2: Draw the initial particle  $\mathbf{z}_0$  from

$$rac{\mu_0(\mathrm{d}oldsymbol{z})}{\mathrm{d}oldsymbol{z}} \propto \exp\left(-L\|oldsymbol{z}\|^2 - rac{\|oldsymbol{x}_0 - e^{-\eta}oldsymbol{z}\|^2}{2(1 - e^{-2\eta})}
ight).$$

- 3: for s = 0 to S 1 do
- Draw a sample  $\tilde{z}_s$  from the Gaussian distribution 4:
- if  $z_{s+1} \notin \mathcal{B}(z_s, r) \cap \mathcal{B}(\mathbf{0}, R)$  then 5:
- $\boldsymbol{z}_{s+1} = \boldsymbol{z}_s;$ 6:
- continue; 7:
- Calculate the accept rate as 8:

$$= \min \left\{ 1, \exp \left( r_g(\boldsymbol{z}_s, \tilde{\boldsymbol{z}}_s) + \frac{\|\tilde{\boldsymbol{z}}_s - \boldsymbol{z}_s + \tau \cdot s_\theta(\boldsymbol{z}_s)\|^2 - \|\boldsymbol{z}_s - \tilde{\boldsymbol{z}}_s + \tau \cdot s_\theta(\tilde{\boldsymbol{z}}_s)\|^2}{4\tau} \right) \right\}$$

Update the particle  $\boldsymbol{z}_{s+1} = \tilde{\boldsymbol{z}}_s$  with probability a, otherwise  $\boldsymbol{z}_{s+1} = \boldsymbol{z}_s$ . 9: 10: return  $\mathbf{z}_S$ ;

1: Input: Returned particle of the previous iteration  $x_0$ , current iteration number k, inner iteration

ribution 
$$\mathcal{N}(\boldsymbol{z}_s - \tau \cdot s_{\theta}(\boldsymbol{z}_s), 2\tau \boldsymbol{I});$$

▷ This condition step is only activated for Projected MALA.

MH Acceptance probability for mitigating the bias

### $r_{q}(\mathbf{z}', \mathbf{z}) \approx \log p_{t}(\mathbf{z}') - \log p_{t}(\mathbf{z})$ denotes the approximation of the energy difference.



### **RTK with MALA: Assumptions**

### **Assumptions:**

- For all  $t \ge 1$ , the score  $\nabla p_t$  is *L*-Lipschitz.
- The second moment of the data distribution  $p_*$  is upper bounded.
- The score and energy difference estimation errors satisfy: for all z and z'

$$\|s_{\theta,t}(\mathbf{z}) - \nabla \log p_t(\mathbf{z})\| \le \epsilon_{\text{score}}$$

$$|r_t(\mathbf{z}',\mathbf{z}) + \log p_t(\mathbf{z}',\mathbf{z})|$$

### MALA algorithm needs both score, i.e., $\nabla \log p$ and the energy difference $\log p(z) - \log p(z')$ .

 $\mathbf{z}') - \log p_t(\mathbf{z}) \leq \epsilon_{\text{energy}}$ 



### **RTK with MALA: Rates**

### **Theorem:** Under the assumptions on the score and energy difference estimation

Suppose the estimation errors of score and energy difference satisfy

$$\epsilon_{\rm score} \leq \frac{\rho \epsilon}{L d^{1/2}} \quad {\rm and} \quad \epsilon_{\rm energy} \leq \frac{\rho \epsilon}{L^2 \cdot \left( d^{1/2} + m_2 + Z \right)},$$

If we implement Alg. 2 with the projected version of MALA with the following hyperparameter settings

$$\eta = \frac{1}{2} \log \frac{2L+1}{2L},$$

and

 $rac{\mu_0(\mathrm{d}oldsymbol{z})}{\mathrm{d}oldsymbol{z}}\propto\mathsf{ex}$ 

it has  $\|\hat{p}_{K\eta} - p_*\|_{TV} \leq \tilde{\mathcal{O}}(\epsilon)$  with an  $\mathcal{O}(L^4 \rho^{-2} \cdot (d + m_2^2)^2 Z^2 \cdot \log(d/\epsilon))$  complexity.

$$K = 4L \cdot \log \frac{(1+L^2)d + \|\nabla f_*(\mathbf{0})\|^2}{\epsilon^2}$$

$$\left\langle \mathsf{p}\left(-L\|\boldsymbol{z}\|^2 - \frac{\|\hat{\boldsymbol{x}}_k - \boldsymbol{e}^{-\eta}\boldsymbol{z}\|^2}{2(1 - \boldsymbol{e}^{-2\eta})}
ight),$$

### **Comparison with Other Methods**

Results	Algorithm	Complexity
Chen et al. (2022)	DDPM (SDE-based)	$ ilde{\mathcal{O}}(L^2 d\epsilon^{-2})$
Chen et al. (2024)	DPOM (ODE-based)	$ ilde{\mathcal{O}}(L^3 d\epsilon^{-2})$
Chen et al. (2024)	DPUM (ODE-based)	$\tilde{\mathcal{O}}(L^2 d^{1/2} \epsilon^{-1})$
Li et al. (2023)	ODE-based sampler	$\tilde{\mathcal{O}}(d^3\epsilon^{-1})$
Ours	RTK-MALA	$O(L^4 d^2 \log(d/\epsilon))$
Ours	RTK-ULD	$\tilde{\mathcal{O}}(L^2 d^{1/2} \epsilon^{-1})$

high-accuracy sampling tasks.

Table. Comparison with prior works for RTK-based methods. The complexity denotes the number of calls for the score estimation to achieve  $\|\hat{p}_{K\eta} - p_*\|_{TV} \leq \tilde{\mathcal{O}}(\epsilon)$ . *d* and  $\epsilon$  mean the dimension and error tolerance. Compared with the state-of-the-art result, RTK-ULD achieves the best dependence for both d and  $\epsilon$ . Though RTK-MALA requires slightly stricter assumptions and worse dimension dependence, a linear convergence w.r.t.  $\epsilon$  makes it suit

### Numerical Experiments: Sampling Gaussian Mixture



### Numerical Experiments: Sampling Gaussian Mixture



### **Autoregressive Diffusion Model**

adapt the result of single diffusion model, we need:

- Each conditional score  $\nabla \log p_t(x | x_1, \dots, x_k)$  should be Lipschitz
- Each target conditional distribution  $p_*(\cdot | x_1, \dots, x_k)$  need to have bounded second moment.
- Each score estimation need to be small.

$$\mathbb{E}_{p_{t}(\cdot|x_{1},...,x_{k})}\left[\|s_{\theta,t}(\mathbf{z};x_{1},...,x_{k}) - \nabla \log p_{t}(\mathbf{z}|x_{1},...,x_{k})\|^{2}\right] \leq \epsilon_{\text{score}}^{2} \quad \forall k = 1,...,K$$

For autoregressive diffusion model, we need to perform K diffusion inference steps. Then, if direct

Do we really need these assumptions for autoregressive diffusion model?

### **Autoregressive Diffusion Model: Assumptions**

In fact, we do not need to make assumptions on all conditional targets, but only require • The target distribution  $p_*$  satisfies  $\|\nabla^2 \log p_*\|_2 \le L$ ,  $\|\nabla \log p_*\|_2 \le \sqrt{L}$ 

- The target distribution  $p_*$  has bounded moment.
- The average of score estimation errors should be small:

$$\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}_{p_t(\cdot|x_1,\dots,x_k)} \left[ \|s_{\theta,t}(\mathbf{z};x_1,\dots,x_k) - \nabla \log p_t(\mathbf{z}|x_1,\dots,x_k)\|^2 \right] \le \epsilon_{\text{score}}^2$$

The assumption is almost as mild as that for standard diffusion model!

# **Autoregressive Diffusion Model: General Theory**

**Theorem:** Under the aforementioned assumptions and apply appropriate learning rate for diffusion sampling, let  $\eta$  be the base stepsize and R be the iteration number in each DM, then considering DDPM solver for diffusion inference, we have

$$\operatorname{KL}\left(p_* \| \hat{p}_*\right) \lesssim 2e^{-2T}L \cdot (m_0 + d) + \left(L^2 R \eta^2 + T \eta\right) \cdot d + \eta m_0 + \eta K R \cdot \epsilon_{\text{score}}^2.$$

Truncation error

- Discretization error is accumulated in dimension:  $d = \sum_{i=1}^{K} d_i$

Huang et al., Capturing Conditional Dependence via Auto-regressive Diffusion Models, ArXiv 2025

Accumulated discretization error Accumulated score estimation error

• The score estimation error is accumulated with an additional K factor (but  $\epsilon_{\text{score}}$  is different).



# **Autoregressive Diffusion Model: General Theory**

**Theorem:** Under the aforementioned assumptions and apply appropriate learning rate for diffusion sampling, let  $\eta$  be the base learning rate and R be the iteration number in each DM, then considering DDPM solver for diffusion inference, we have

$$\mathrm{KL}\left(p_* \| \hat{p}_*\right) \lesssim 2e^{-2T} L \cdot (m_0 + d) + (L^2 R \eta^2 + T \eta) \cdot d + \eta m_0 + \eta K R \cdot \epsilon_{\mathrm{score}}^2.$$

Accumulated discretization error

- assumptions.)
  - For DDPM, the total complexity will be  $\tilde{O}(Kd\epsilon^{-2})$ .
  - For RTK-ULD, the total complexity will be  $\tilde{O}(Kd^{1/2}\epsilon^{-1})$ .

Autoregressive diffusion model requires longer inference time!

• The discretization error can be improved using RTK-ULD or RTK-MALA (need stronger

### **Autoregressive Diffusion Model: Conditional Dependency**

**Lemma:** Under the aforementioned assumptions and apply appropriate learning rate for diffusion sampling, considering DDPM solver, we have

$$\begin{split} & \operatorname{KL}\left(p_{*,k+1|[1:k]}(\cdot|\boldsymbol{x}_{[1:k]})\big\|\hat{p}_{*,k+1|[1:k]}(\cdot|\boldsymbol{x}_{[1:k]})\right) \\ & \lesssim e^{-2T} \cdot \left(2Ld_{k+1} + \mathbb{E}_{p_{*,k+1|[1:k]}(\cdot|\boldsymbol{x}_{[1:k]})}\left[\|\mathbf{y}\|^{2}\right]\right) + \eta \cdot \sum_{r=0}^{R-1} \tilde{L}_{k+1,r}(\boldsymbol{\theta}|\boldsymbol{x}_{[1:k]}) + d_{k+1}L^{2}R\eta^{2} \\ & + d_{k+1}T\eta + \eta \mathbb{E}_{p_{*,k+1|[1:k]}(\cdot|\boldsymbol{x}_{[1:k]})}\left[\|\mathbf{y}\|^{2}\right]. \end{split}$$

- For any conditional distribution, AR DM can well capture it with appropriate choice of step size.
- This implies that the dependency between tokens  $x_i$  and  $x_j$  can be well reflected in the generated data  $x = [x_1, \dots, x_K]$

# Single Diffusion Model: Possible Failure

**Lemma 4.4.** Consider random vectors  $\mathbf{y} \in \mathbb{R}^{d_{k+1}}$  and  $\mathbf{x} \in \mathbb{R}^{d_1+d_2+\cdots+d_k}$ . For any error threshold  $\varepsilon \in (0, 1/2]$  and for any  $M \in \mathbb{R}$ , there exists a pair of Gaussian probability densities  $(p_*(\mathbf{y}, \mathbf{x}), \hat{p}_*(\mathbf{y}, \mathbf{x})), \text{ such that } \mathrm{KL}(p_*(\mathbf{y}, \mathbf{x}) \| \hat{p}_*(\mathbf{y}, \mathbf{x})) \leq \varepsilon, \text{ while } \mathrm{KL}(p_*(\mathbf{y}|\mathbf{x}) \| \hat{p}_*(\mathbf{y}|\mathbf{x})) > \varepsilon$  $M^2 \cdot \| \boldsymbol{x}_{(1:d_{k+1})} \|^2.$ 

- learning the conditional dependency between different parts in the image.
- A simple case is  $x = [x_1, x_2], x_1 \sim N(0)$ 
  - conditioned).
  - Th

• Single diffusion model can well capture the joint distribution, while may be failed in

(0,1), 
$$x_2 = x_1 + N(0,\epsilon^2)$$

• The joint distribution x has very large condition number (covariance matrix is kind of ill-

2



# **Autoregressive Diffusion Model: Conditional Dependency**



- $\bullet$ capture the inter-patch dependencies, and achieves lower training loss.
- This aligns with the theoretical results on conditional dependency.

When the dependency appear among different patches of the data, AR DM can better



**'**O







# **Numerical Experiments: Task 1**



- This aligns with the theoretical results in general setting.

• When the dependency does not appear among different patches of the data, AR DM can not capture the inter-patch dependencies, and achieves substantially larger training loss.



### Takeaway

- AR diffusion model and Standard diffusion model
  - We developed an RTK framework, which can be used to develop a family of diffusion inference algorithms.
  - We developed the sampling theory for AR diffusion model, which shows that: AR diffusion model behaves worse than standard DM, in general settings. AR diffusion model can better capture inter-patch dependencies than
- - standard DM.