

Evolution and Application of Diffusion Models (DMs) in Reinforcement Learning (RL)

Networked Intelligence for Comprehensive Efficiency (NICE) Lab
College of Information Science and Electronic Engineering
Zhejiang University
<http://nice.rongpeng.info/>



Sep. 18, 2025

Content



1 Research Background & Preliminaries

- Deep Generative Models
- Diffusion Models
- Different Roles of DMs in RL

2 DMs in Single-Agent RL

- Single-Agent: DMs as Policy
- Single-Agent: DMs as Planner

3 DMs in Multi-Agent RL

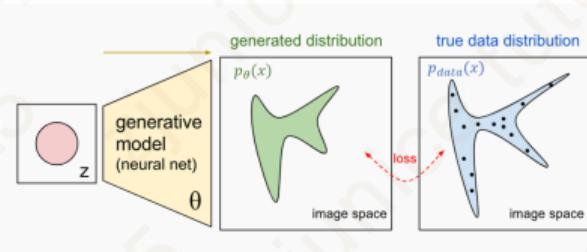
- Multi-Agent: DMs as Policy
- Multi-Agent: DMs as Planner
- Multi-Agent: DMs as A Unified Framework

4 Conclusion & Future Prospect

Research Background & Preliminaries



Deep Generative Models¹



■ Maximizing Log-Likelihood

$$\begin{aligned}\theta^* &= \arg \min_{\theta} D_{\text{KL}}(p_{\text{data}}(x) \| p_{\theta}(x)) \\ &= \arg \min_{\theta} \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log p_{\text{data}}(x) - \log p_{\theta}(x)] \\ &= \arg \max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log p_{\theta}(x)] \\ &\approx \arg \max_{\theta} \sum_{i=1}^N \log p_{\theta}(x_i)\end{aligned}$$

- ▶ $p(x) = \int p(x, z) dz$: the marginalization integral over the latent variable z is intractable in complex models
- ▶ $p(x) = \frac{p(x, z)}{q(z|x)}$: the posterior distribution is difficult to obtain directly



■ Optimizing Evidence Lower Bound (ELBO)

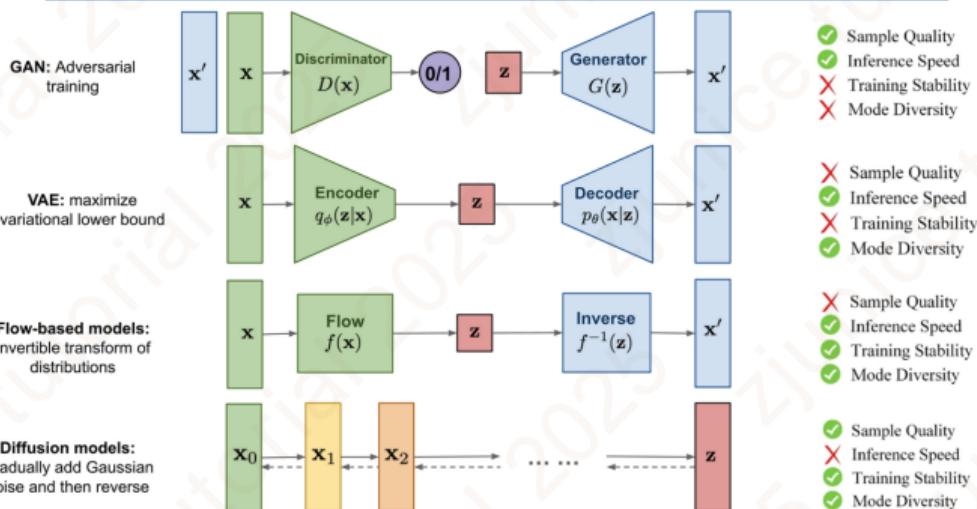
$$\begin{aligned}\log p(x) &= \int q_{\phi}(z|x) \log p(x) dz \quad (\text{Bring Evidence into Integral}) \\ &= \mathbb{E}_{q_{\phi}(z|x)} [\log p(x)] \quad (\text{Definition of Expectation}) \\ &= \mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{p(x, z)}{p(z|x)} \right] \quad (\text{Apply Bayes Rule}) \\ &= \mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{p(x, z)}{p(z|x) q_{\phi}(z|x)} \right] \\ &= \mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{p(x, z)}{q_{\phi}(z|x)} \right] + \mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{q_{\phi}(z|x)}{p(z|x)} \right] \\ &= \mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{p(x, z)}{q_{\phi}(z|x)} \right] + D_{\text{KL}}(q_{\phi}(z|x) \| p(z|x)) \\ &\geq \boxed{\mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{p(x, z)}{q_{\phi}(z|x)} \right]} \quad (\text{KL Divergence} \geq 0)\end{aligned}$$

¹L. Ruthotto and E. Haber, "An introduction to deep generative modeling," *GAMM-Mitteilungen*, vol. 44, no. 2, e202100008, 2021.



Deep Generative Models²

$$\begin{aligned}
 \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{p(x, z)}{q_\phi(z|x)} \right] &= \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{p_\theta(x|z)p(z)}{q_\phi(z|x)} \right] && \text{(Chain Rule of Probability)} \\
 &= \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)] + \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{p(z)}{q_\phi(z|x)} \right] && \text{(Split the Expectation)} \\
 &= \underbrace{\mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)]}_{\text{reconstruction term}} - \underbrace{D_{KL}(q_\phi(z|x) \parallel p(z))}_{\text{prior matching term}} && \text{(Definition of KL Divergence)}
 \end{aligned}$$



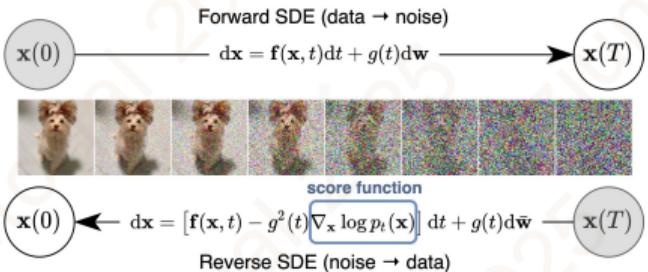
²Z. Ahmad et al., “Understanding gans: Fundamentals, variants, training challenges, applications, and open problems,” *Multimedia Tools and Applications*, vol. 84, no. 12, pp. 10 347–10 423, 2025 .

DDIM

Diffusion Models^{3,4}



Score-Based



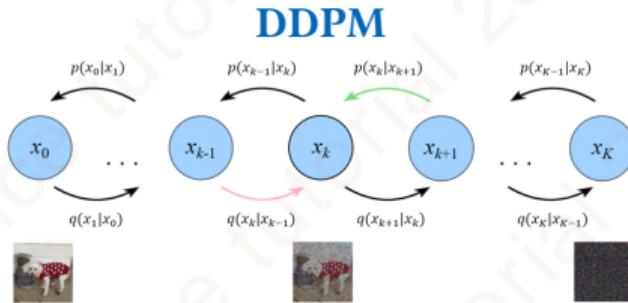
- Forward SDE:

$$dx_t = f(x_t, t) dt + g(t) d\omega_t$$

- Reverse SDE:

$$dx_t = [f(x_t, t) - g^2(t) \boxed{\nabla \log p_t(x_t)}] dt + g(t) d\bar{\omega}_t$$

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{t, x_0, x_t} [\lambda(t) \|s_\theta(x_t, t) - \nabla \log p(x_t | x_0)\|_2^2]$$



- Forward diffusion process:

$$q(x_k | x_{k-1}) \sim \mathcal{N}(\sqrt{\alpha_k} x_{k-1}, (1 - \alpha_k) I)$$

$$q(x_k | x_0) \sim \mathcal{N}(\sqrt{\alpha_k} x_0, (1 - \bar{\alpha}_k) I)$$

- Reverse denoising process:

$$q(x_{k-1} | x_k, x_0) = \frac{q(x_k | x_{k-1}, x_0) q(x_{k-1} | x_0)}{q(x_k | x_0)}$$

$$\sim \mathcal{N}\left(\frac{1}{\sqrt{\alpha_k}} x_k - \frac{1 - \alpha_k}{\sqrt{\alpha_k} \sqrt{1 - \bar{\alpha}_k}} \boxed{\epsilon_0}, \sigma_k^2 I\right)$$

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{t, x_0, x_k} [\lambda(k) \|\epsilon_k \boxed{\epsilon_0}\|_2^2]$$

³Y. Song et al., "Score-based generative modeling through stochastic differential equations," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, Virtual Edition, 2021.

⁴J. Ho et al., "Denoising diffusion probabilistic models," in *Proc. Adv. Neural Inf. Proces. Syst. (NeurIPS)*, Virtual Edition, 2020.



Diffusion Models⁵

■ Feature 1: Reparameterization

Rewrite a random variable as a **deterministic function of a noise variable**. This enables gradient-based optimization of non-stochastic terms.

$$z \sim \mathcal{N}(z; \mu_\theta, \sigma_\theta^2 I)$$



$$z = \mu_\theta + \sigma_\theta \odot \epsilon, \quad \epsilon \sim \mathcal{N}(0, I)$$

■ Feature 2: x_k can be directly represented by x_0

Forward Process:

$$\begin{aligned} x_k &= \sqrt{\alpha_k} x_{k-1} + \sqrt{1 - \alpha_k} \epsilon_{k-1} \\ &= \sqrt{\alpha_k \alpha_{k-1}} x_{k-2} + \sqrt{1 - \alpha_k \alpha_{k-1}} \epsilon_{k-2} \\ &= \dots \\ &= \sqrt{\bar{\alpha}_k} x_0 + \sqrt{1 - \bar{\alpha}_k} \epsilon_0 \quad \bar{\alpha}_k = \prod_{i=1}^k \alpha_i \end{aligned}$$

Reconstruct:

$$x_0 = \frac{x_k - \sqrt{1 - \bar{\alpha}_k} \epsilon_0}{\sqrt{\bar{\alpha}_k}}$$

⁵J. Ho et al., "Denoising diffusion probabilistic models," in *Proc. Adv. Neural Inf. Proces. Syst. (NeurIPS)*, Virtual Edition, 2020.

Denoising Diffusion Implicit Model (DDIM)⁶

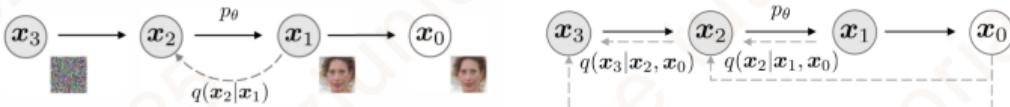


Figure 1: Graphical models for diffusion (left) and non-Markovian (right) inference models.

■ Non-Markovian sampling process

$$\begin{aligned} q_\sigma(x_{k-1} | x_k, x_0) \\ = \mathcal{N}\left(\sqrt{\bar{\alpha}_{k-1}}x_0 + \sqrt{1 - \bar{\alpha}_{k-1} - \sigma_k^2} \frac{x_k - \sqrt{\bar{\alpha}_k}x_0}{\sqrt{1 - \bar{\alpha}_k}}, \sigma_k^2 I\right) \end{aligned}$$

The mean is chosen to ensure that

$$q_\sigma(x_k | x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_k}x_0, (1 - \bar{\alpha}_k)I)$$

- When $\sigma \rightarrow 0$, x_{k-1} is uniquely determined by x_0 and x_k .
- When $\sigma_k = \sqrt{\frac{1 - \bar{\alpha}_{k-1}}{1 - \bar{\alpha}_k}}(1 - \alpha_k)$, the generalized sampler is equivalent to DDPM.

■ Generalized generative process

$$\text{Reconstruct: } \hat{x}_0 = \frac{x_k - \sqrt{1 - \bar{\alpha}_k} \epsilon_\theta(x_k, k)}{\sqrt{\bar{\alpha}_k}}$$

$$\begin{aligned} x_{k-1} &= \sqrt{\bar{\alpha}_{k-1}} \frac{x_k - \sqrt{1 - \bar{\alpha}_k} \epsilon_\theta(x_k, k)}{\sqrt{\bar{\alpha}_k}} \\ &+ \sqrt{1 - \bar{\alpha}_{k-1} - \sigma_k^2} \epsilon_\theta(x_k, k) + \sigma_k \epsilon_k \end{aligned}$$

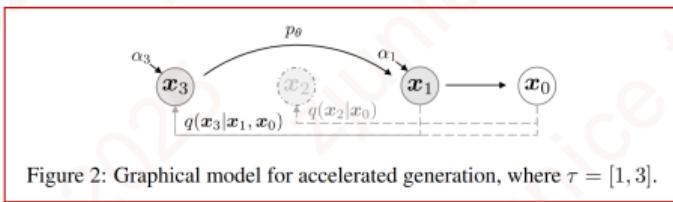


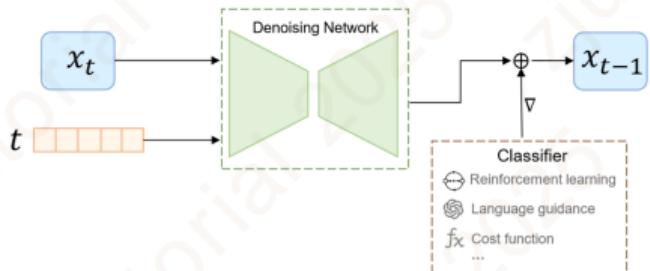
Figure 2: Graphical model for accelerated generation, where $\tau = [1, 3]$.

⁶J. Song et al., “Denoising diffusion implicit models,” in *Proc. Int. Conf. Learn. Represent. (ICLR)*, Virtual Edition, 2021.



Conditional Diffusion Model^{7,8}

Classifier Guidance

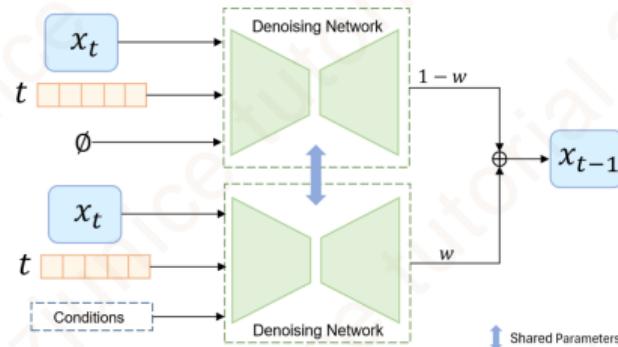


$$\nabla \log p_t(\mathbf{x}_t | \mathbf{y}) = \nabla \log p_t(\mathbf{x}_t) + \nabla \log p_t(\mathbf{y} | \mathbf{x}_t)$$

$$\mathbf{x}_{t-1} \sim \mathcal{N}\left(\mu_{\theta}(\mathbf{x}_t, t) + w \boxed{\nabla_{\mathbf{x}_t} \log p_{\phi}(\mathbf{y} | \mathbf{x}_t, t)}, \sigma_t^2 \mathbf{I}\right)$$

- Stronger control
- Higher sample quality

Classifier-free Guidance



$$\epsilon_t = w \epsilon_{\theta}(\mathbf{x}_t, t, \mathbf{y}) + (1 - w) \epsilon_{\theta}(\mathbf{x}_t, t, \emptyset)$$

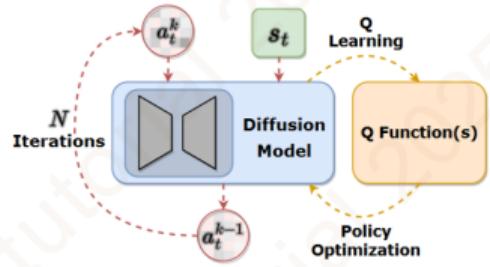
- No need for extra classifier
- More stable and flexible
- Preserves sample diversity

⁷P. Dhariwal and A. Nichol, "Diffusion models beat gans on image synthesis," in *Proc. Adv. Neural Inf. Proces. Syst. (NeurIPS)*, Virtual Edition, 2021.

⁸J. Ho and T. Salimans, "Classifier-free diffusion guidance," *arXiv preprint arXiv:2207.12598*, 2022 .

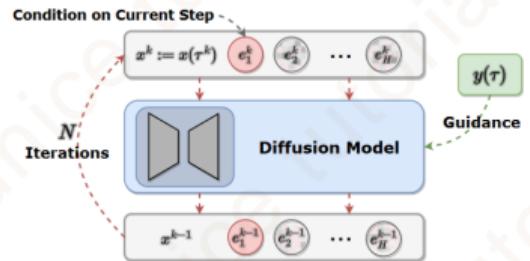


Different Roles of DMs in RL⁹



DM as policy \longleftrightarrow **MFRL**

- Expressive multimodal policy parameterization
- Improved exploration and diversity
- Robustness to perturbations through iterative denoising



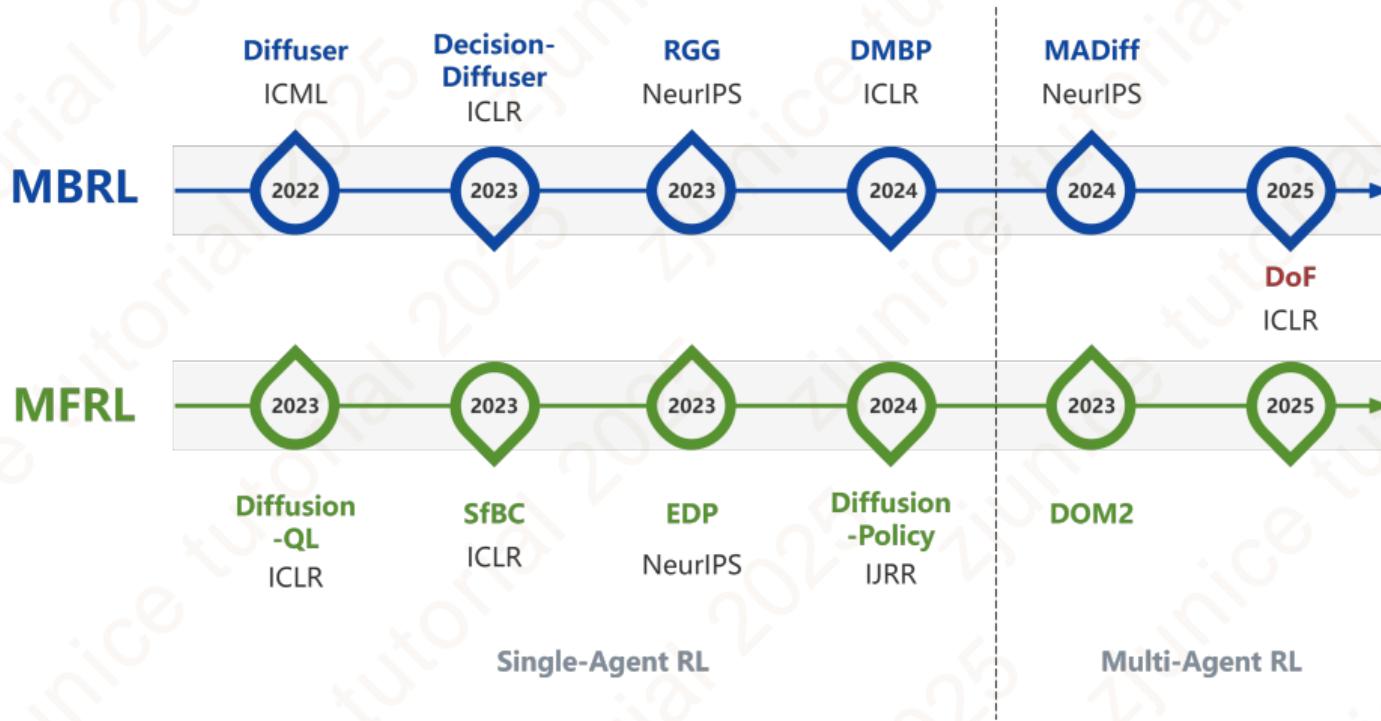
DM as planner \longleftrightarrow **MBRL**

- Global optimization from a long-term perspective
- Flexibility in multimodal trajectory generation
- Enhanced robustness through iterative refinement

⁹Z. Zhu et al., "Diffusion models for reinforcement learning: A survey," *arXiv preprint arXiv:2311.01223*, 2023 .

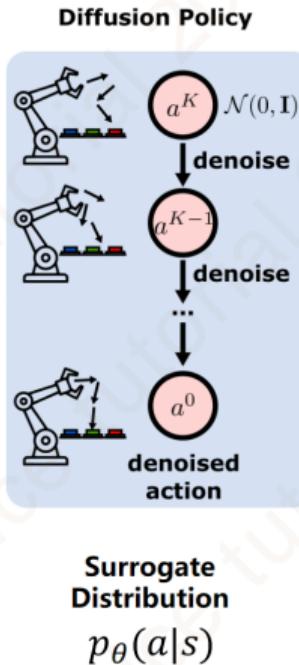


Algorithmic Timeline of DMs in RL



DMs in Single-Agent RL

Diffusion-QL¹⁰



- Behavior-cloning loss

$$\mathcal{L}_d(\theta) = \mathbb{E}_{k, \epsilon, (s, a)} \left[\left\| \epsilon - \epsilon_\theta \left(\sqrt{\bar{\alpha}_k} \boxed{a} + \sqrt{1 - \bar{\alpha}_k} \epsilon, \boxed{s}, k \right) \right\|^2 \right]$$

- Q-value function loss

$$\mathbb{E}_{(s_t, a_t, s_{t+1}) \sim \mathcal{D}, a_{t+1}^0 \sim \pi_{\theta'}} \left[\left\| \left(r(s_t, a_t) + \gamma \min_{i=1,2} Q_{\phi'_i}(s_{t+1}, a_{t+1}^0) \right) - Q_{\phi_i}(s_t, a_t) \right\|^2 \right]$$

$$\begin{aligned} \pi &= \arg \min_{\pi_\theta} \mathcal{L}(\theta) = \mathcal{L}_d(\theta) + \mathcal{L}_q(\theta) \\ &= \underbrace{\mathcal{L}_d(\theta)}_{\text{policy regularization}} - \underbrace{\alpha \cdot \mathbb{E}_{s \sim D, a^0 \sim \pi_\theta} [Q_\phi(s, a^0)]}_{\text{policy improvement}} \end{aligned}$$

$\frac{\partial Q}{\partial a}$ is backpropagated through the whole diffusion chain

¹⁰Z. Wang et al., "Diffusion policies as an expressive policy class for offline reinforcement learning," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, Kigali, Rwanda, 2023.



Efficient Diffusion Policy (EDP)¹¹

■ Action reconstruction

Forward: $\mathbf{x}_k = \sqrt{\bar{\alpha}_k} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_k} \boldsymbol{\epsilon}_0$

$$\hat{a}_0 = \frac{1}{\sqrt{\bar{\alpha}_k}} a_k - \frac{\sqrt{1 - \bar{\alpha}_k}}{\sqrt{\bar{\alpha}_k}} \epsilon_\theta(a_k, k, s)$$

■ Policy improvement

$$L_\pi(\theta) = -\mathbb{E}_{s \sim \mathcal{D}, \hat{a}_0} [Q_\phi(s, \hat{a}_0)]$$

Direct policy optimization:

$$\nabla_\theta L_\pi(\theta) = -\frac{\partial Q_\phi(s, a)}{\partial a} \frac{\partial a}{\partial \theta}$$

$\frac{\partial a}{\partial \theta}$ is tractable

Likelihood-based policy optimization:

$$\max_{\theta} \mathbb{E}_{(s,a) \sim D} [f(Q_\phi(s, a)) \log \pi_\theta(a | s)]$$

$$\Leftrightarrow \mathbb{E}_{k,\varepsilon,(a,s)} \left[\frac{\beta^k f(Q_\phi(s, a))}{2 \alpha^k (1 - \bar{\alpha}^{k-1})} \| \varepsilon - \varepsilon_\theta(a^k, k; s) \|^2 \right]$$

$$\Leftrightarrow \mathbb{E}_{k,\varepsilon,(a,s)} [f(Q_\phi(s, a)) \| a - \hat{a}^0 \|^2]$$

$\log \pi_\theta(a | s)$ is intractable

Optimizing the ELBO in DDPM

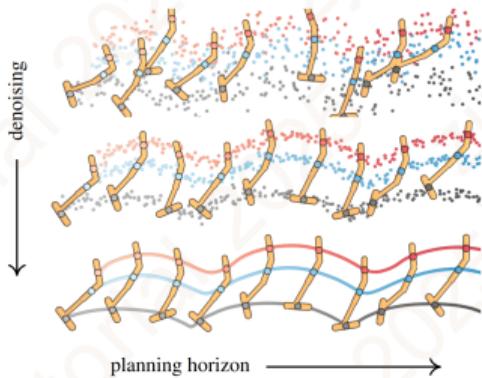
Policy approximation
 $\hat{\pi}_\theta(a | s) \triangleq \mathcal{N}(\hat{a}^0, I)$



¹¹B. Kang et al., "Efficient diffusion policies for offline reinforcement learning," in Proc. Adv. Neural Inf. Proces. Syst. (NeurIPS), New Orleans, Louisiana, USA, 2023.



Diffuser¹²

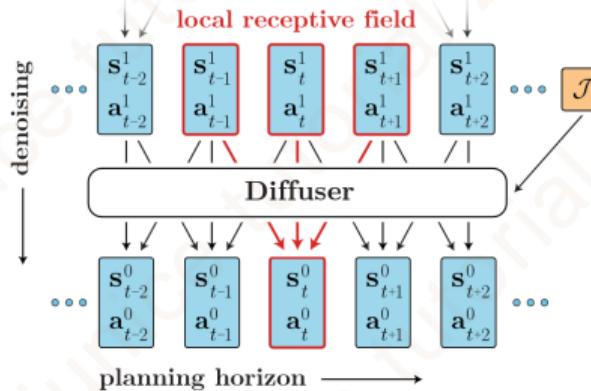


■ Diffusing over state-action trajectory

$$\mathbf{x}_k = \begin{bmatrix} s_0 & s_1 & \dots & s_H \\ a_0 & a_1 & \dots & a_H \end{bmatrix}$$

■ Guiding with classifier guidance

$$\begin{aligned} p_{\theta}(\mathbf{x}_{k-1} | \mathbf{x}_k, \mathbf{y}) &\approx \mathcal{N}(\mathbf{x}_{k-1}; \boldsymbol{\mu}_{\theta} + \Sigma \boxed{\mathbf{g}}, \Sigma) \\ \mathbf{g} &= \nabla_{\tau} \log p(\mathbf{y} | \tau) \Big|_{\tau=\boldsymbol{\mu}_{\theta}} \\ &= \nabla \mathcal{J}(\boldsymbol{\mu}_{\theta}) \end{aligned}$$



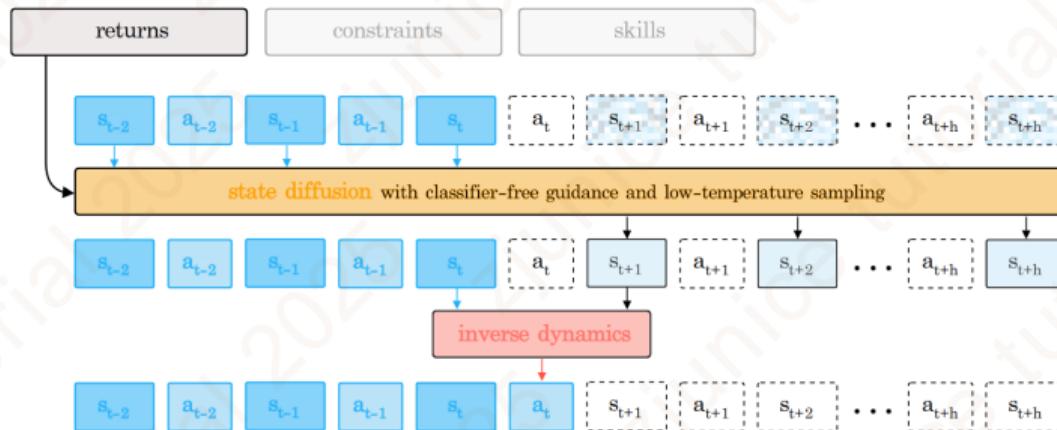
Limitations

- ▷ Actions are harder to predict and model
- ▷ Requires estimating a Q-function

¹²M. Janner et al., "Planning with diffusion for flexible behavior synthesis," in Proc. Int. Conf. Mach. Learn. (ICML), Baltimore, Maryland, USA, 2022.



Decision-Diffuser¹³



- Diffusing over state trajectory

$$\mathbf{x}_k(\tau) = (s_t, s_{t+1}, \dots, s_{t+H-1})_k$$

- Acting with inverse-dynamics

$$a_t = f_\phi(s_t, s_{t+1})$$

- Guiding with classifier-free guidance

$$\hat{\epsilon} = \epsilon_\theta(\mathbf{x}_k(\tau), \emptyset, k)$$

$$+ \omega \left(\epsilon_\theta(\mathbf{x}_k(\tau), \mathbf{y}(\tau), k) - \epsilon_\theta(\mathbf{x}_k(\tau), \emptyset, k) \right)$$

¹³A. Ajay et al., “Is conditional generative modeling all you need for decision-making?” In Proc. Int. Conf. Learn. Represent. (ICLR), Kigali, Rwanda, 2023.

DMs in Multi-Agent RL



DMs in Multi-Agent RL



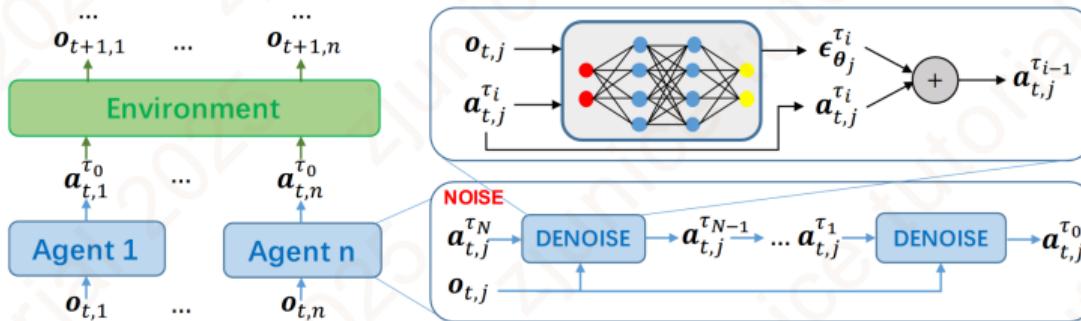
- MADiff-C

- DOM2

- MAD2RL
- MADiff-D
- DoF



DOM²¹⁴



■ Diffusion-QL-based policy update

- DTDE: agent j learns and runs its own policy θ_j

$$\mathcal{L}_{bc}(\theta_j) = \mathbb{E}_{(o_{t,j}, a_{t,j}^0), \epsilon, k} \left[\left\| \epsilon - \epsilon_{\theta_j} \left(a_{t,j}^k, o_{t,j}, k \right) \right\|_2^2 \right]$$

- Sample based on 1-order DPM-solver (DDIM)

$$a_{t,j}^{k-1} = \frac{\bar{\alpha}_{k-1}}{\bar{\alpha}_k} a_{t,j}^k - \sigma_k \left(\frac{\bar{\alpha}_k \sigma_{k-1}}{\bar{\alpha}_{k-1} \sigma_k} - 1 \right) \epsilon_{\theta_j}$$

■ CQL-based Q-value update

$$\begin{aligned} \mathcal{L}(\phi_j) = & \mathbb{E}_{(o_j, a_j) \sim \mathcal{D}_j} [(Q_{\phi_j}(o_j, a_j) - y_j)^2] \\ & + \zeta \mathbb{E}_{(o_j, a_j) \sim \mathcal{D}_j} [\log \sum_{\tilde{a}_j} \exp(Q_{\phi_j}(o_j, \tilde{a}_j)) - Q_{\phi_j}(o_j, a_j)] \end{aligned}$$

suppress the unseen state-action pairs
encourage the state-action pairs from the dataset

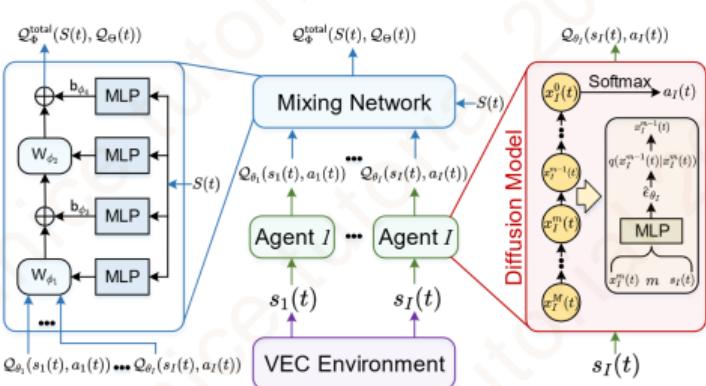
¹⁴Z. Li et al., "Beyond conservatism: Diffusion policies in offline multi-agent reinforcement learning," *arXiv preprint arXiv:2307.01472*, 2023.

MAD2RL¹⁵



Incorporate DMs to determine the optimal DNN partitioning and task offloading decisions in Vehicular Edge Computing (VEC).

- DM-based action making



$$Q_{\theta_i}^u(s_i(t)) = \frac{e^{x_i^{0,u}(t)}}{\sum_{\nu} e^{x_i^{0,\nu}(t)}}, \quad a_i(t) = \arg \max_u Q_{\theta_i}^u(s_i(t))$$

- QMIX-based policy improvement

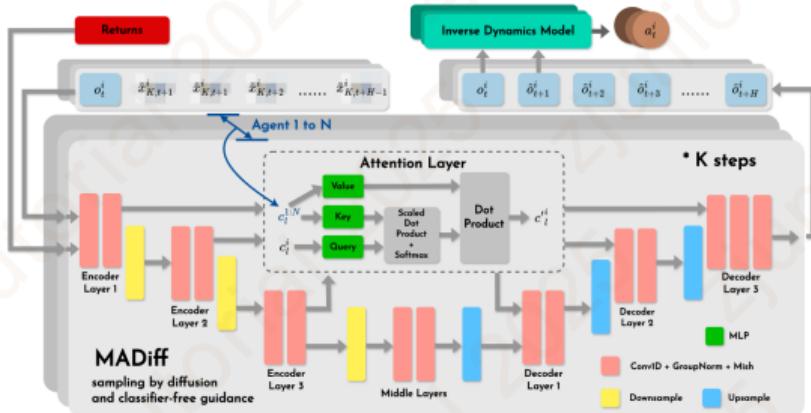
$$Q_{\Theta}(t) = \left\{ Q_{\theta_1}(s_1(t), a_1(t)), \dots, Q_{\theta_N}(s_N(t), a_N(t)) \right\}$$

$$Q_{\Phi}^{\text{total}} = \text{MixingNET}(S(t), Q_{\Theta}(t))$$

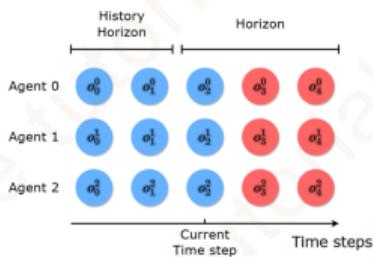
¹⁵Z. Liu et al., "Dnn partitioning, task offloading, and resource allocation in dynamic vehicular networks: A lyapunov-guided diffusion-based reinforcement learning approach," *IEEE Trans. Mob. Comput.*, pp. 1–17, 2024 .



MADiff¹⁶



■ Execution



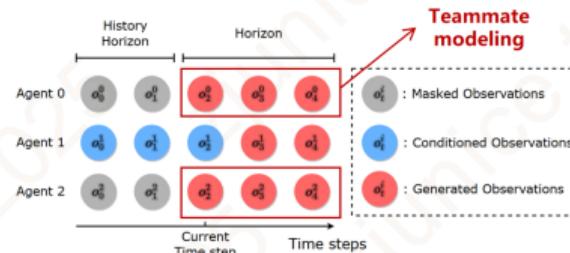
(a) MADIFF in centralized control.

■ Information interchange

Multi-head attention on the skip-connected feature from symmetric l -th encoder layer of agent i : \mathbf{c}_l^i .

■ Centralized training

$$\begin{aligned} \mathcal{L}(\theta, \phi) := & \sum_i \mathbb{E}_{(s^i, a^i, s'^i) \in \mathcal{D}} [\|a^i - I_\phi^i(o^i, o'^i)\|^2] & \text{individual inverse dynamics model} \\ & + \mathbb{E}_{k, \tau_0 \in \mathcal{D}, \beta} [\|\epsilon - \epsilon_\theta(\hat{\tau}_k, (1 - \beta)y(\tau_0) + \beta\emptyset, k)\|^2] & \text{unified diffusion noise model} \end{aligned}$$

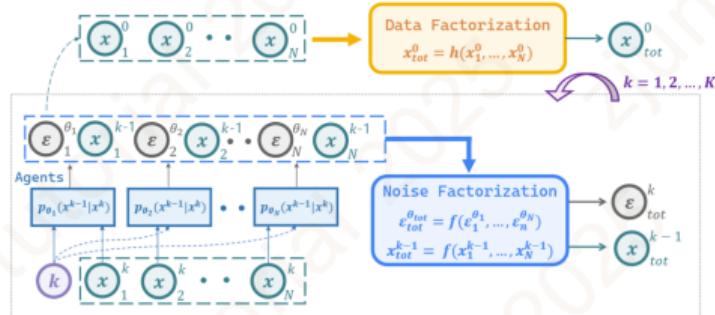


(b) MADIFF in decentralized execution.

¹⁶Z. Zhu et al., "Madiff: Offline multi-agent learning with diffusion models," in Proc. Adv. Neural Inf. Proces. Syst. (NeurIPS), Vancouver, CANADA, 2024.



DoF¹⁷



Diffusion factorization

Factorize the joint noise and joint data as

$$\begin{aligned}\epsilon_{tot}^k &= f(\epsilon_1^k, \dots, \epsilon_N^k) \quad 0 \leq k \leq K \\ \epsilon_{tot}^{\theta_{tot}}(x_{tot}^k, k) &= f(\epsilon_1^{\theta_1}(x_1^k, k), \dots, \epsilon_N^{\theta_N}(x_N^k, k)) \quad 0 \leq k \leq K \\ x_{tot}^k &= f(x_1^k, \dots, x_N^k) \quad 1 \leq k \leq K \\ x_{tot}^0 &= h(x_1^0, \dots, x_N^0)\end{aligned}$$

Individual-Global-Identically-Distributed (IGD)

$$\arg \max_{\boldsymbol{u}} Q_{tot}(\boldsymbol{\tau}, \boldsymbol{u}) = \begin{pmatrix} \arg \max_{u_1} Q_1(\tau_1, u_1) \\ \vdots \\ \arg \max_{u_N} Q_N(\tau_N, u_N) \end{pmatrix}$$

► Individual-Global-Identically-Distributed (IGD)

Assume individual marginals

$$p_{\theta_i}(x_i^0) = \int p_{\theta_i}(x_i^{0:K}) dx_i^{1:K}, \text{ such that}$$

$$\prod_{i=1}^N p_{\theta_i}(x_i^0) = p_{\theta_{tot}}(x_{tot}^0), \quad \theta_i \subset \theta_{tot}.$$

Then the generated samples $\{x_i^0\}_{i=1}^N$ share the same distribution as x_{tot}^0 .

¹⁷C. Li et al., "Dof: A diffusion factorization framework for offline multi-agent reinforcement learning," in Proc. Int. Conf. Learn. Represent. (ICLR), Singapore, 2025.



DoF

► Concat

$$x_{\text{tot}}[(i-1) \times d : i \times d] = x_i$$

► WConcat

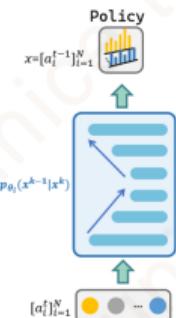
$$x_{\text{tot}}[(i-1) \times d : i \times d] = k_i x_i$$

► Atten

$$x_{\text{tot}}[(i-1) \times d : i \times d] = \sum_{j=1}^N w_i^j x_j$$

► QMIX

$$x_{\text{tot}} = h(x_1, x_2, \dots, x_n)$$



■ DoF as policy

Follow the design of [Diffusion-QL](#):

$$\begin{aligned}\mathcal{L} &= \mathcal{L}_{\text{diff}}(\theta) + \mathcal{L}_{\text{pg}}(\theta) \\ &= \mathbb{E} \left[\|\epsilon - \epsilon_{\theta}^{\text{tot}}\|^2 \right] - \alpha \mathbb{E} \left[Q_{\phi}(s, u^0) \right]\end{aligned}$$

Algorithm 1 Centralized Training

```

1: repeat
2:    $\mathbf{x}_{\text{tot}}^0 \sim q(\mathbf{x}_{\text{tot}})$  (Sample global data)
3:    $k \sim \text{Uniform}(\{1, \dots, K\})$  (Diffusion step)
4:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \in \mathbb{R}^{d \times N}$  (Sample global noise)
5:    $\mathbf{x}_{\text{tot}}^k = \sqrt{\alpha^k} \mathbf{x}_{\text{tot}}^{k-1} + \sqrt{1 - \alpha^k} \epsilon$ 
6:    $\mathbf{x}_i^k = \mathbf{x}_{\text{tot}}^k[(i-1) \times d : i \times d], i \in [1, \dots, N]$ 
7:    $\epsilon_{\text{tot}} = f(\epsilon_{\theta_1}^1(\mathbf{x}_1^k, k), \dots, \epsilon_{\theta_N}^N(\mathbf{x}_N^k, k))$ 
8:   Take gradient descent step on:
      
$$\nabla_{\theta} \|\epsilon - \epsilon_{\text{tot}}\|^2$$

9: until convergence

```

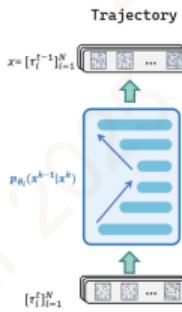
Algorithm 2 Decentralized Execution

```

1:  $\mathbf{x}_i^K \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  (Initialize for each agent i)
2: for  $k = K, \dots, 1$  do
3:    $\epsilon_{\theta}^k(\mathbf{x}_i^k, k)$  (Noise prediction by each agent i)
4:   Update state for each agent  $i$ :
      
$$\mathbf{x}_i^{k-1} = \frac{1}{\sqrt{\alpha_k}} \left( \mathbf{x}_i^k - \frac{1 - \alpha_k}{\sqrt{1 - \alpha_k}} \epsilon_{\theta}^k(\mathbf{x}_i^k, k) \right) + \sigma_k \mathbf{z},$$

      where  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $k > 1$ , else  $\mathbf{z} = \mathbf{0}$ .
5: end for
6: return  $\mathbf{x}_i^0$  (Final trajectory or action for each agent i)

```



■ DoF as planner

Follow the design of [Decision-Diffuser](#):

$$\begin{aligned}\mathcal{L} &= \mathbb{E} \left[\|\epsilon - \epsilon_{\theta}^{\text{tot}}(\mathbf{x}_k, (1-\beta)\mathbf{y} + \beta\emptyset, k)\|^2 \right] \\ &\quad + \mathbb{E} \left[\|\mathbf{u} - D_{\phi}(\mathbf{o}, \mathbf{o}')\|^2 \right]\end{aligned}$$

Conclusion & Future Prospect



Conclusion

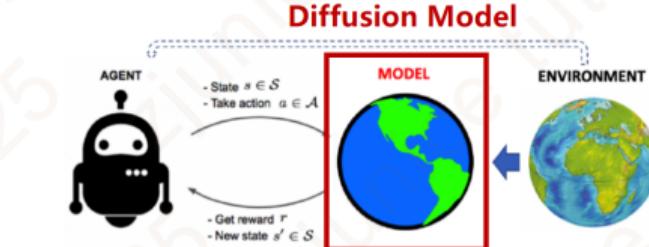
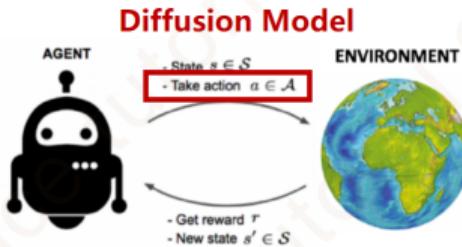
I have talked about

■ Fundamentals of Diffusion Models

- ▶ Superior generative modeling power over other deep generative models
- ▶ Mathematical foundations of diffusion models
- ▶ Sampling acceleration techniques for diffusion models
- ▶ Conditional generative modeling with diffusion models

■ Applications of Diffusion Models in Reinforcement Learning

- ▶ Diffusion models as policy
- ▶ Diffusion models as planner



- ▶ Diffusion models as a unified framework

Future Prospect



■ Fast Inference and Online Adaptation of DMs

DMs require multi-step iterative denoising, which leads to low inference efficiency and limited real-time scalability.

DDIM; DPM-Solver; DPM-Solver++; Flow; Consistency Model ...

■ Optimize the Architecture of the Noise Model

Most current noise models adopt simple backbones like MLPs or U-Nets, which struggle to capture the complex interactions among agents.

Attention mechanism (enhance inter-agent interaction);

LSTM (better historical insight) ...

■ Generalization to Multi-Task Scenarios

DMs can be post-trained to adapt to varying tasks and environments, enabling better generalization in multi-task settings.

Importance sampling; beam search; conditional guidance; RL-based fine-tuning ...

Thank you

Networked Intelligence for Comprehensive Efficiency (NICE) Lab
College of Information Science and Electronic Engineering

Zhejiang University
<https://nice.rongpeng.info>