



澳門大學  
UNIVERSIDADE DE MACAU  
UNIVERSITY OF MACAU

# Diffusion Policy: Visuomotor Policy Learning via Action Diffusion

GONG Ying

[gong.ying@connect.umac.mo](mailto:gong.ying@connect.umac.mo)

2024.1.12

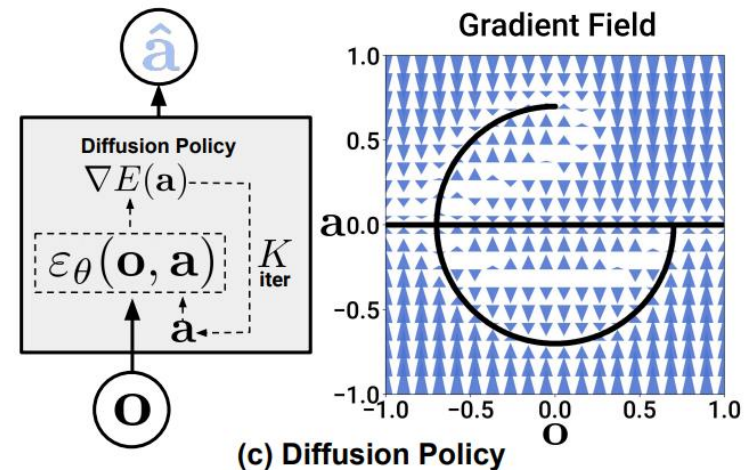
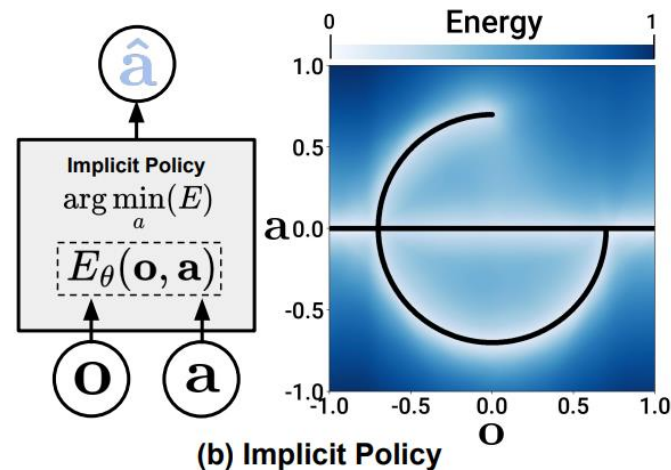
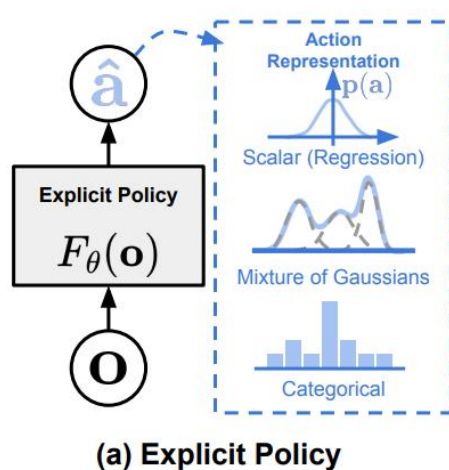
# Contents

1. Introduction
2. Diffusion Policy Formulation
3. Key Design Decisions
4. Diffusion Policy Overview



# 1. Introduction

- Policy learning from demonstration can be formulated as the supervised regression task of learning to **map observations to actions**.
- Prior work attempts to explore different action representations from explicit to implicit to better capture multi-modal distributions.
- This work introduces a new form of robot visuomotor policy that generates behavior via a "conditional denoising diffusion process on robot action space", **Diffusion Policy**.



# 1. Introduction

- **Diffusion policy** infers the action-score gradient, conditioned on visual observations, for K denoising iterations.
  - I. Expressing multimodal action distributions.
  - II. High-dimensional output space.
  - III. Stable training.
- **Contributions:**
  - I. Closed-loop action sequences.
  - II. Visual conditioning.
  - III. Time-series diffusion transformer.
- Consistent performance boost across all benchmarks with an average improvement of **46.9%**. (12 tasks from 4 benchmarks)



## 2. Diffusion Policy Formulation

Visuomotor robot policies are formulated as **Denoising Diffusion Probabilistic Models (DDPMs)**.

- **Denoising Diffusion Probabilistic Models**

Starting from  $x^K$  sampled from Gaussian noise, the DDPM performs  $K$  iterations of denoising to produce a series of actions with decreasing levels of noise,  $x^k, x^{k-1} \dots x^0$ , until a desired noise-free output  $x^0$  is formed.

$$x^{k-1} = \alpha(x^k - \gamma \varepsilon_\theta(x^k, k) + N(0, \sigma^2 I)),$$

where  $\varepsilon_\theta$  is the noise prediction network and  $N(0, \sigma^2 I)$  is Gaussian noise added. It can be interpreted as a noisy gradient descent step:

$$x' = x - \gamma \nabla E(x).$$

Where  $\varepsilon_\theta(x, k)$  predicts the gradient field  $\nabla E(x)$  and  $\gamma$  is the learning rate. An  $\alpha$  slightly smaller than 1 improves stability.



## 2. Diffusion Policy Formulation

- **DDPM Training**

The training process starts by randomly drawing unmodified examples,  $\mathbf{x}^0$ , from the dataset. For each sample, we randomly select a denoising iteration  $k$  and then sample a random noise  $\varepsilon^k$  with appropriate variance for iteration  $k$ . The noise prediction network is asked to predict the noise from the data sample with noise added.

$$L = \text{MSE} \left( \varepsilon^k, \varepsilon_{\theta}(\mathbf{x}^0 + \varepsilon^k, k) \right),$$

Minimizing the loss function  $L$  also minimizes the variational lower bound of the KL-divergence between the data distribution  $p(\mathbf{x}^0)$  and the distribution of samples drawn from the DDPM  $q(\mathbf{x}^0)$ .



## 2. Diffusion Policy Formulation

- **Diffusion for Visuomotor Policy Learning**

DDPMs are typically used for image generation ( $x$  is an image), we use a DDPM to learn robot visuomotor policies. It requires 2 modifications:

- I. Change the output  $x$  to represent robot actions:

**Closed-loop action-sequence prediction**

- II. Make the denoising processes conditioned on input observation  $O_t$ :

**Visual observation conditioning**



## 2. Diffusion Policy Formulation

- **Diffusion for Visuomotor Policy Learning**

- I. **Closed-loop action-sequence prediction**

To encourage temporal consistency and smoothness in long-horizon planning while allowing prompt reactions to unexpected observations, the action-sequence prediction produced by a diffusion model is integrated with **receding horizon control**:

At time step  $t$  the policy takes the latest  $T_O$  steps of observation data  $\mathbf{O}_T$  as input and predicts  $T_P$  steps of actions, of which  $T_a$  steps of actions are executed on the robot **without re-planning**. Here we define  $T_O$  as the observation horizon,  $T_P$  as the action prediction horizon and  $T_a$  as the action execution horizon.

$$\boxed{o_1, o_2, o_3}, \boxed{o_4}, o_5, \dots, \boxed{o_{t-2}, o_{t-1}, o_t}$$

$T_O = 3$



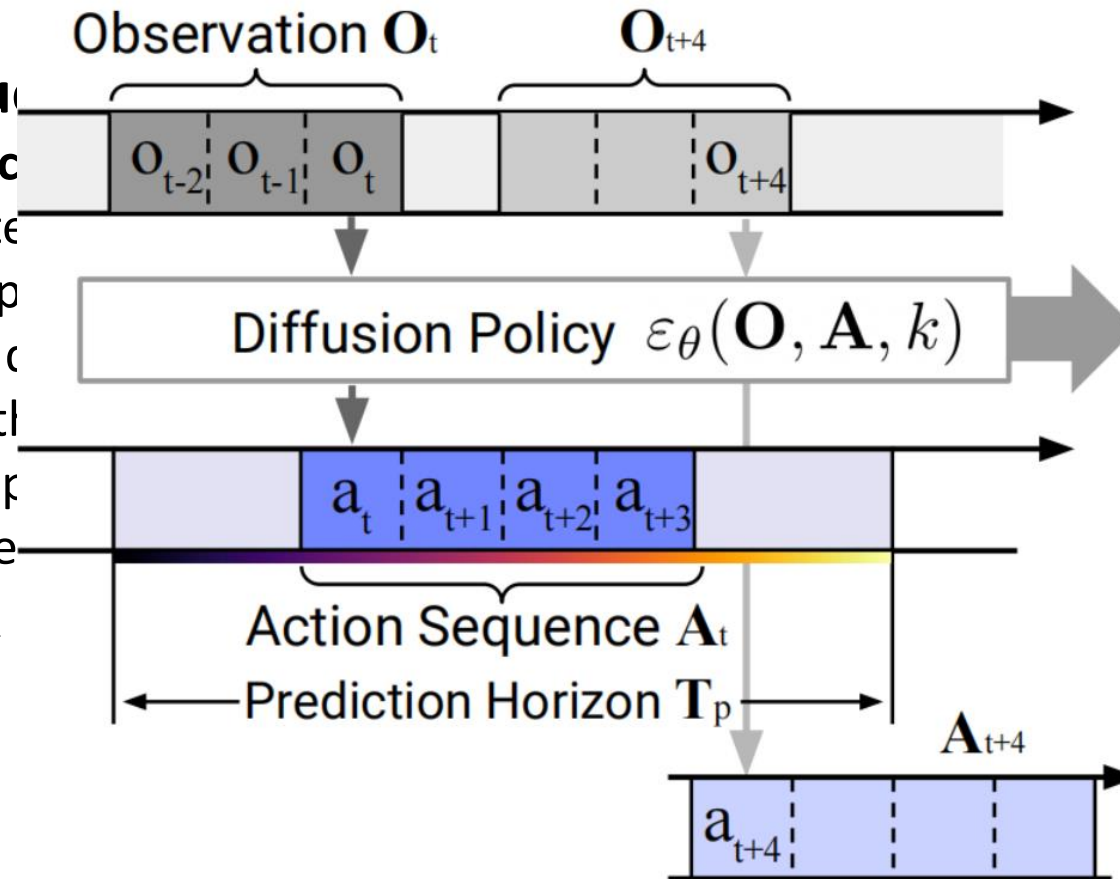


## 2. Diffusion Policy Formulation

- Diffusion for Visuomotor Control

- I. Closed-loop action prediction

To encourage the policy to learn long-term planning while allowing prompt re-planning when the environment changes, we produce by a diffusion policy. At time step  $t$  the policy predicts  $T_p$  steps ahead without re-planning. He uses a prediction horizon  $T_p$  and a planning horizon  $T_a$ .



1 planning while  
sequence prediction  
control:  
ta  $O_T$  as input and  
ed on the robot without  
e action prediction

## 2. Diffusion Policy Formulation

- Diffusion for Visuomotor Policy Learning

- II. Visual observation conditioning

We use a DDPM to approximate the conditional distribution  $p(\mathbf{A}_t|\mathbf{O}_t)$  instead of the joint distribution  $p(\mathbf{A}_t, \mathbf{O}_t)$ . It allows the model to predict actions conditioned on observations **without inferring** future states, speeding up the diffusion process and improving the accuracy of generated actions. To capture  $p(\mathbf{A}_t|\mathbf{O}_t)$ , we modify  $\mathbf{x}^{k-1} = \alpha(\mathbf{x}^k - \gamma \varepsilon_\theta(\mathbf{x}^k, k) + N(0, \sigma^2 I))$  to:

$$\mathbf{A}_t^{k-1} = \alpha(\mathbf{A}_t^k - \gamma \varepsilon_\theta(\mathbf{O}_t, \mathbf{A}_t^k, k) + N(0, \sigma^2 I)),$$

And the training loss is modified from  $L = \text{MSE}(\varepsilon^k, \varepsilon_\theta(x^0 + \varepsilon^k, k))$  to:

$$L = \text{MSE}(\varepsilon^k, \varepsilon_\theta(\mathbf{O}_t, \mathbf{A}_t^0 + \varepsilon^k, k)),$$

The exclusion of  $\mathbf{O}_t$  from the output of the denoising process significantly improves inference speed and better accommodates real-time control.



# 3. Key Design Decisions

- **Network Architecture Options**

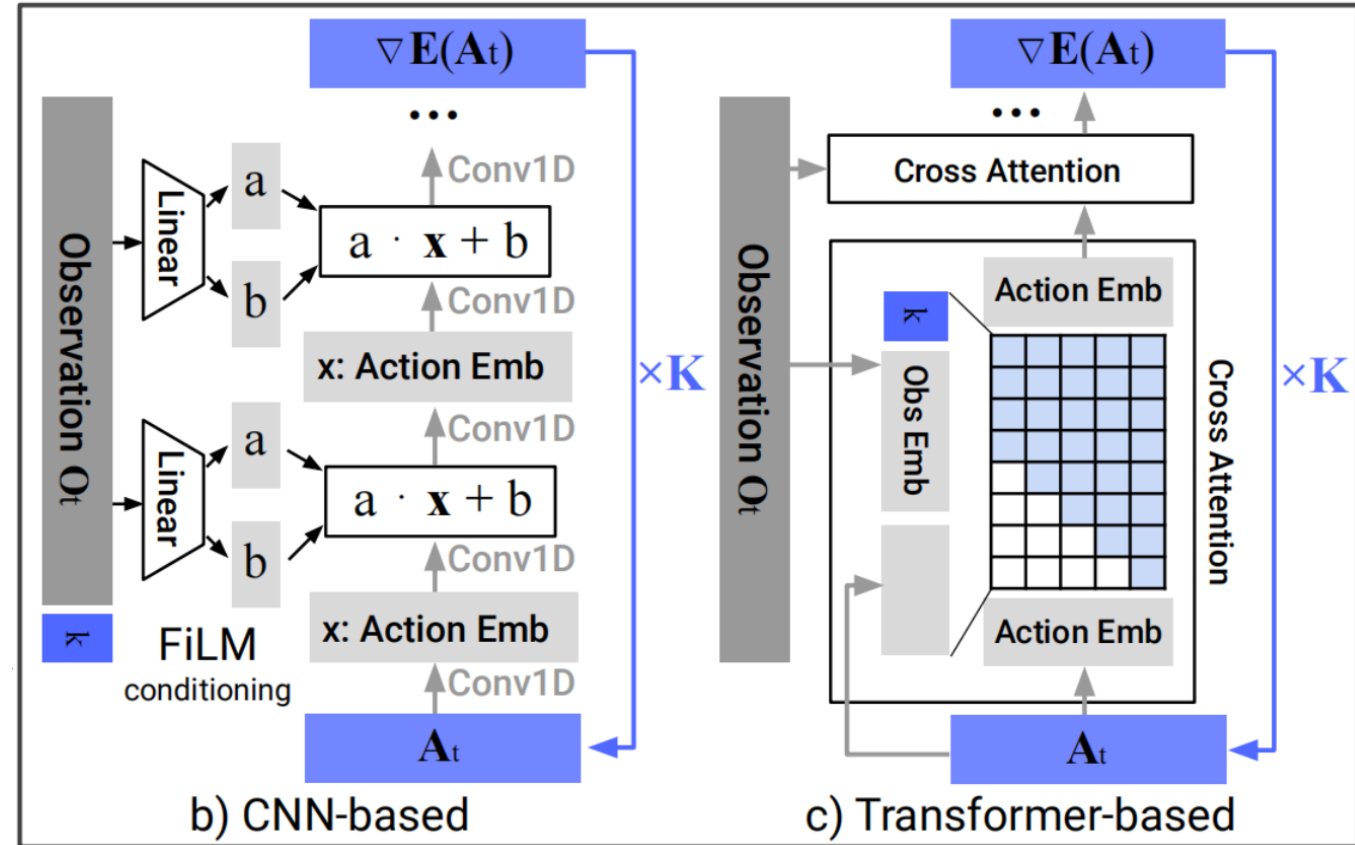
Choice of NN architectures for  $\epsilon_\theta$ .

- I. **CNN-based diffusion policy**

Needs less tuning, but performs poorly when the desired action sequence changes quickly and sharply.

- II. **Time-series diffusion transformer**

More sensitive to hyperparameters, but performs well if the task is complex or action changes often.



# 3. Key Design Decisions

- **Visual Encoder**

- The visual encoder maps the raw **image sequence** into a latent embedding  $\mathbf{O}_t$  and is trained end-to-end with the diffusion policy.
- **ResNet-18** is used as the encoder with the following modifications:
  - I. Replace the global average pooling with a spatial softmax pooling to maintain spatial information.
  - II. Replace BatchNorm with GroupNorm for stable training.

- **Noise Schedule**

It's empirically found that the Square Cosine Schedule proposed in iDDPM works best for the tasks.

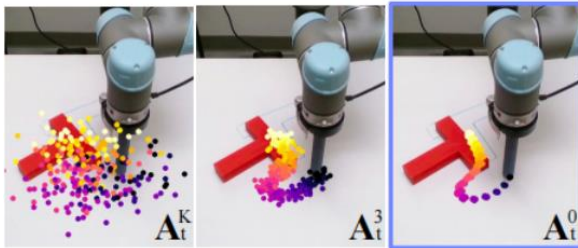
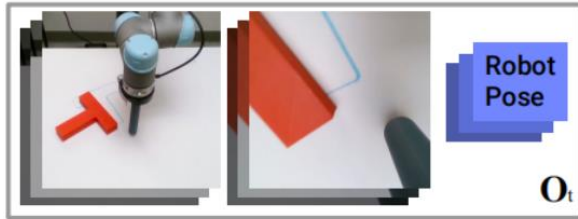
- **Accelerating Inference for Real-time Control**

Denoising Diffusion Implicit Models (DDIM) approach.

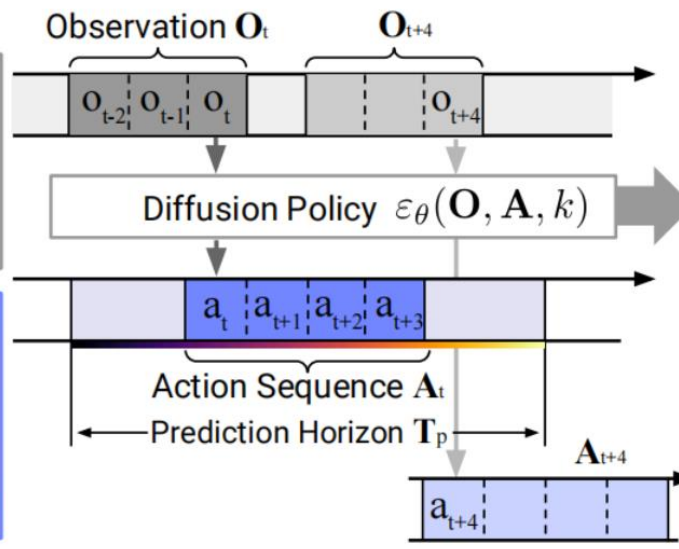


# 4. Diffusion Policy Overview

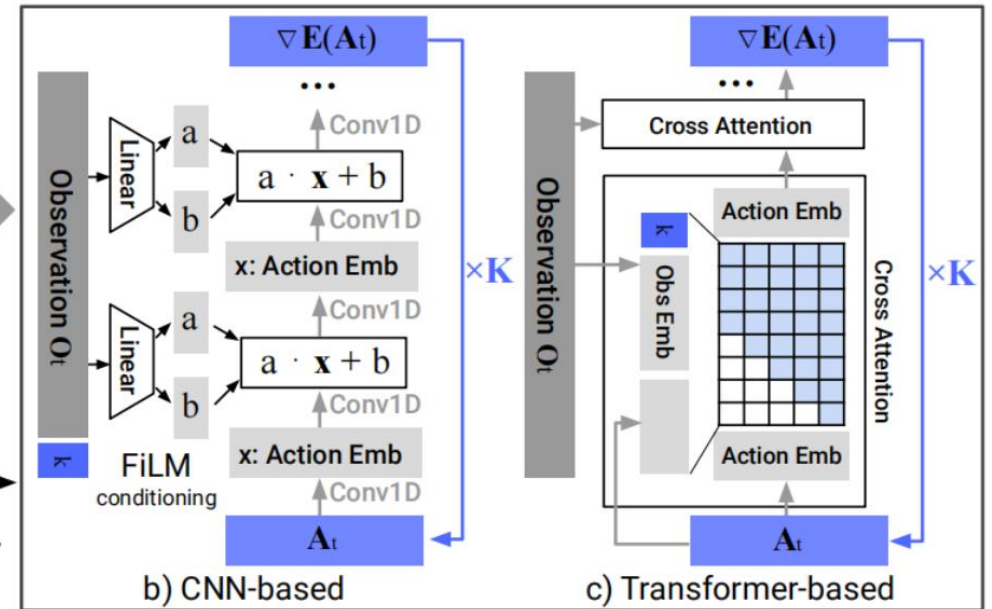
**Input: Image Observation Sequence**



**Output: Action Sequence**



a) Diffusion Policy General Formulation



Thank you.

