

DR

DeepRob

[Group 8] Lecture 6

Diffusion Models and Policy Learning

by *Nirshal Chandra Sekar and Mohit Yadav*

University of Minnesota



Overview

- Motivation behind diffusion models
- Denoising Diffusion Probabilistic Models
 - Diffusion Process and Reverse Diffusion Process
 - Training and Sampling a Diffusion Model
 - Results from a Diffusion Model that we trained
- Conditioned v/s Unconditioned Diffusion
- Policy Learning
- Toy Problem
- Diffusion Policy on Real Robots
 - Advantages
 - Network Architect
 - Evaluation & Results



Motivation Behind Diffusion Model

Text Input ~ "Generate an image of me delivering a presentation on diffusion models."

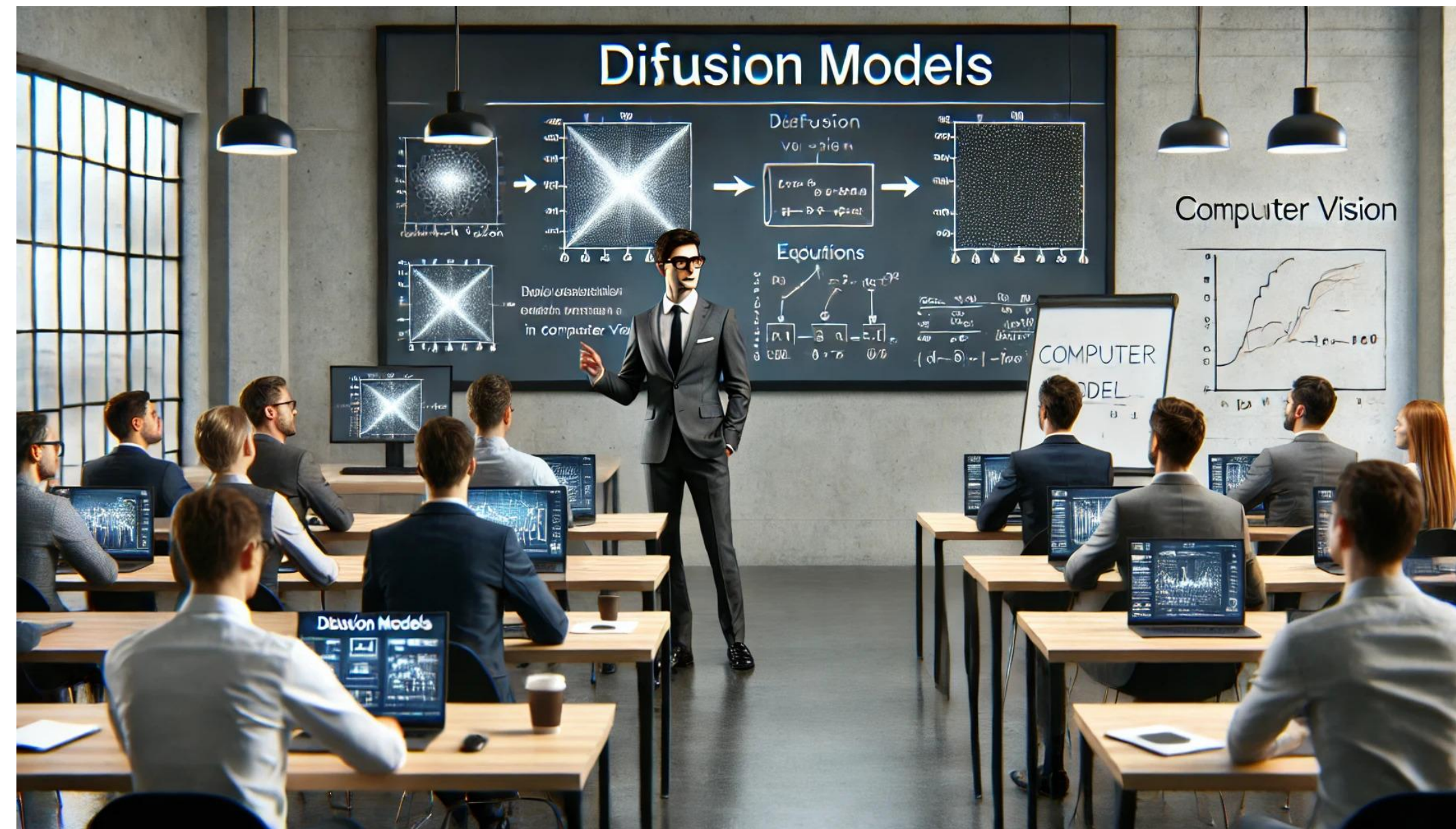


Image Output ~ Generated using ChatGPT - 4o

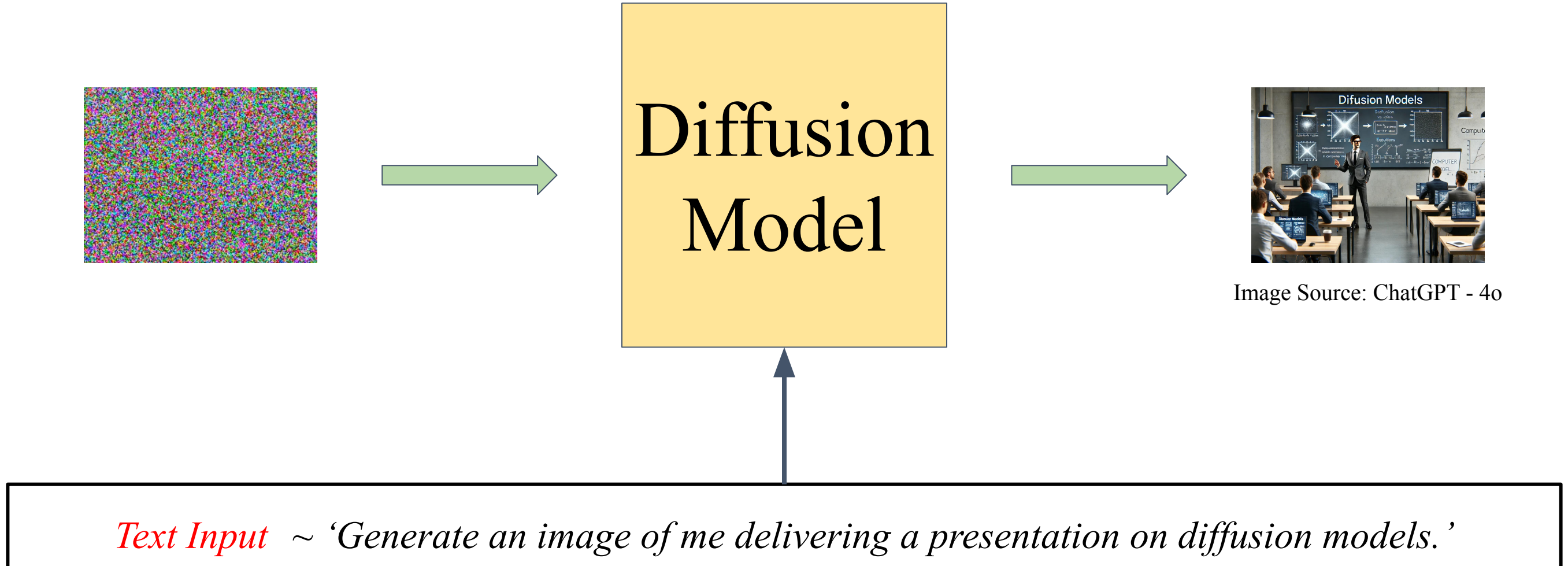
Motivation Behind Diffusion Model

Text Input ~ “Generate an image where the attendees are getting bored during my presentation on diffusion policies.”

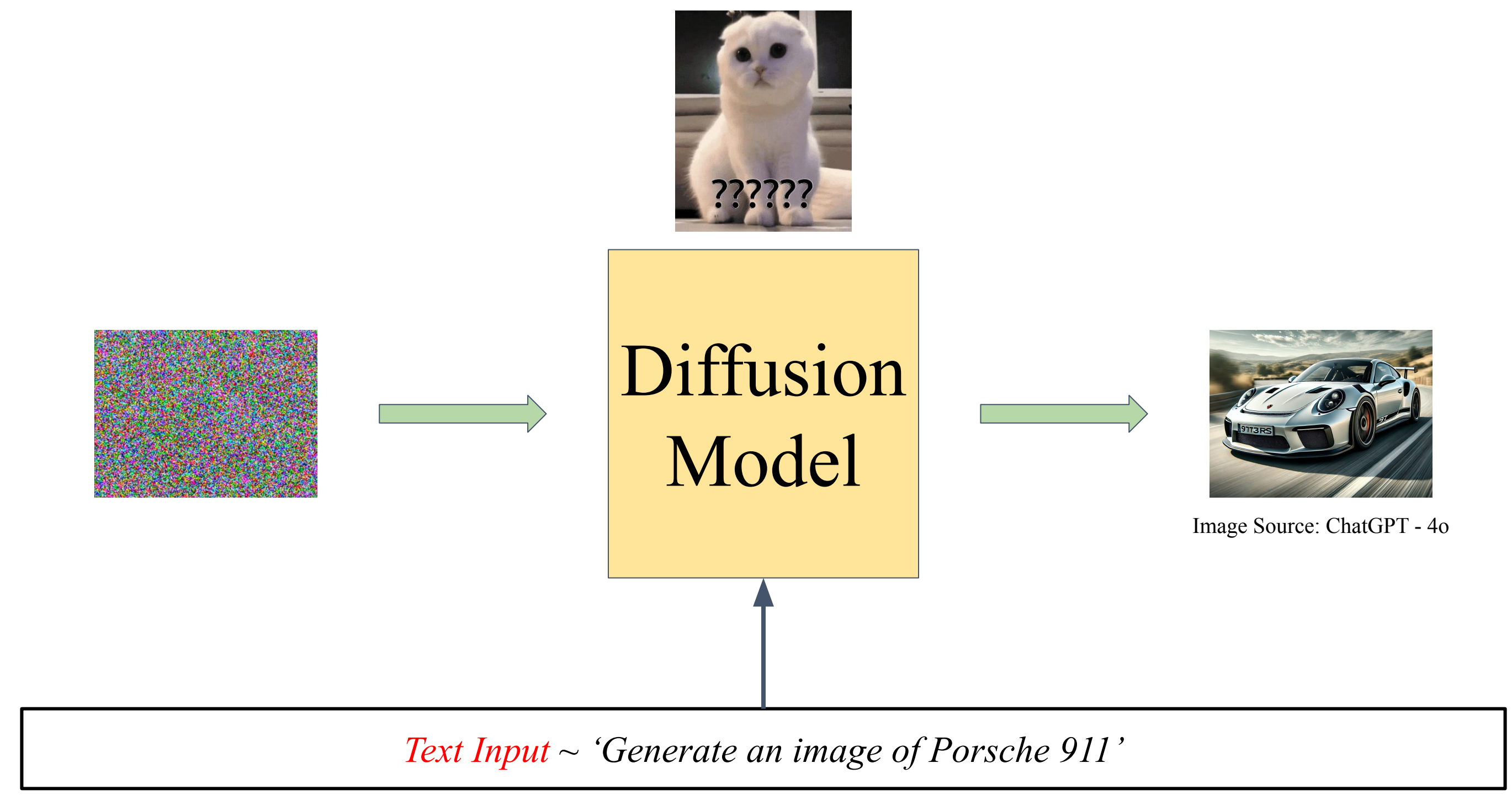


Image Output ~ Generated using ChatGPT - 4o

Motivation Behind Diffusion Model



Motivation Behind Diffusion Model





Overview

- Motivation behind diffusion models ✓
- Denoising Diffusion Probabilistic Models
 - Diffusion Process and Reverse Diffusion Process
 - Training and Sampling a Diffusion Model
 - Results from a Diffusion Model that we trained
- Conditioned v/s Unconditioned Diffusion
- Policy Learning
- Toy Problem
- Diffusion Policy on Real Robots
 - Advantages
 - Network Architect
 - Evaluation & Results





Denoising Diffusion Probabilistic Models

Denoising Diffusion Probabilistic Models

Jonathan Ho
UC Berkeley
jonathanho@berkeley.edu

Ajay Jain
UC Berkeley
ajayj@berkeley.edu

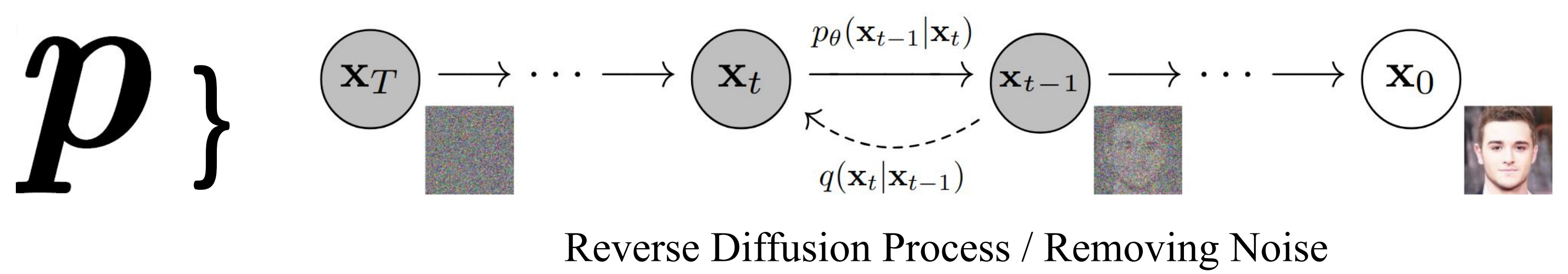
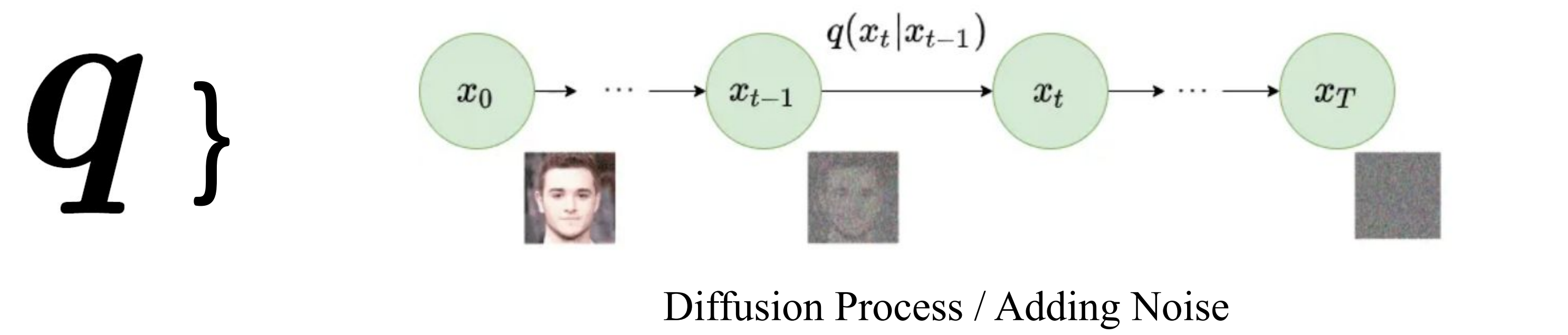
Pieter Abbeel
UC Berkeley
pabbeel@cs.berkeley.edu

Abstract

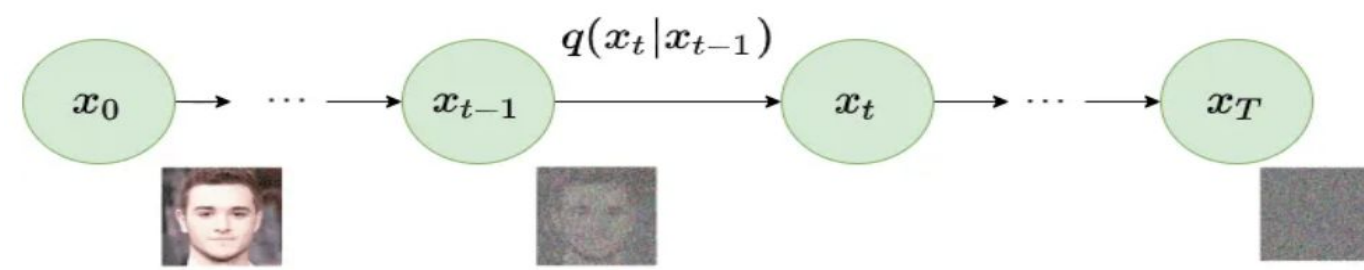
We present high quality image synthesis results using diffusion probabilistic models, a class of latent variable models inspired by considerations from nonequilibrium thermodynamics. Our best results are obtained by training on a weighted variational bound designed according to a novel connection between diffusion probabilistic models and denoising score matching with Langevin dynamics, and our models naturally admit a progressive lossy decompression scheme that can be interpreted as a generalization of autoregressive decoding. On the unconditional CIFAR10 dataset, we obtain an Inception score of 9.46 and a state-of-the-art FID score of 3.17. On 256x256 LSUN, we obtain sample quality similar to ProgressiveGAN. Our implementation is available at <https://github.com/hojonathanho/diffusion>.



Denoising Diffusion Probabilistic Models



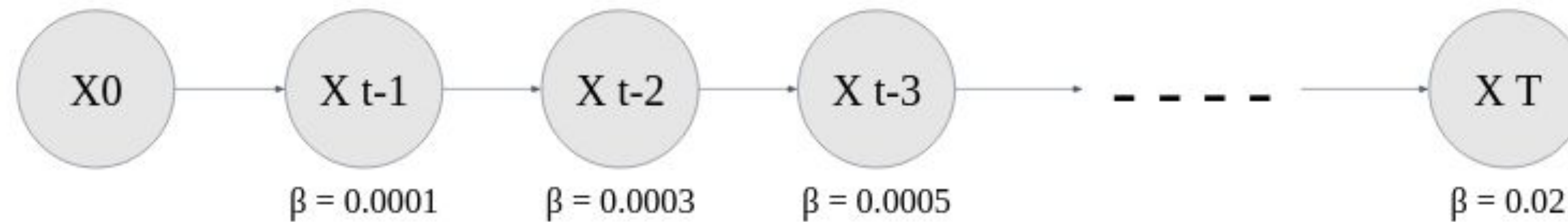
Diffusion Process / Adding Noise



Diffusion Process / Adding Noise

Joint probability of the states x_1, x_2, \dots, x_T in a Markov chain, given the initial state x_0 .

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) := \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}), \quad q(\mathbf{x}_t|\mathbf{x}_{t-1}) := \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$$

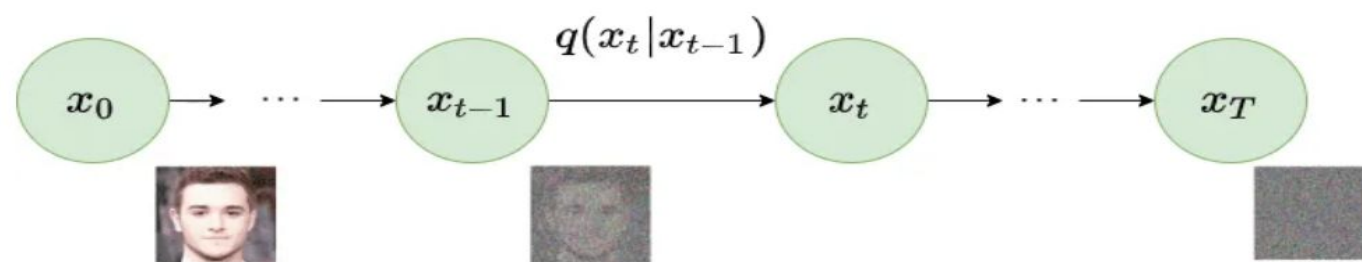


$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \sqrt{1 - \beta_t} \cdot \mathbf{x}_{t-1} + \sqrt{\beta_t} \cdot \text{noise}$$



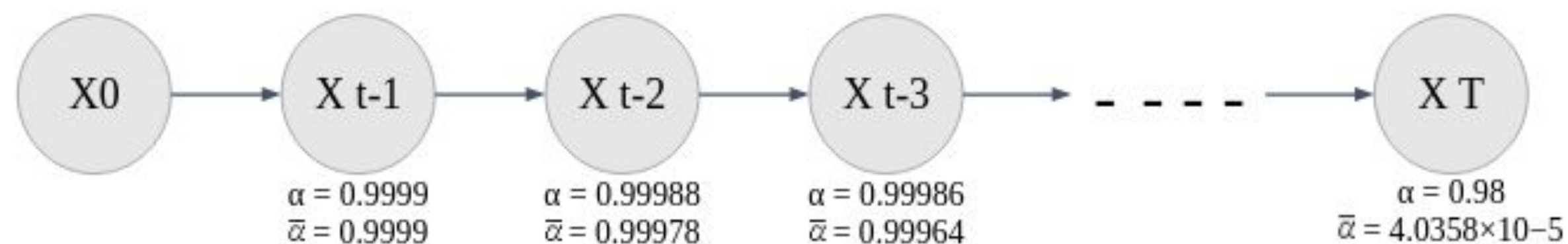


Diffusion Process / Adding Noise



Diffusion Process / Adding Noise

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) := \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}), \quad q(\mathbf{x}_t|\mathbf{x}_{t-1}) := \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$$



$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \sqrt{1 - \beta_t} \cdot \mathbf{x}_{t-1} + \sqrt{\beta_t} \cdot \text{noise}$$

$$q(\mathbf{x}_t|\mathbf{x}_0) = \sqrt{\bar{\alpha}_t} \cdot \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \cdot \text{noise}$$

For more information please refer the paper

Where,
 $\alpha = 1 - \beta$
 $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$





Diffusion Process / Adding Noise

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \sqrt{1 - \beta_t} \cdot x_0 + \sqrt{\beta_t} \cdot \text{noise}$$

- $\beta_1 = 0.0001, \beta_2 = 0.0003, \beta_3 = 0.0005$
- noise = 5
- input = 10

$$q(X_1|X_0) = (\sqrt{1-0.0001} \cdot 10) + (\sqrt{0.0001} \cdot 5) = 10.0494$$

$$q(X_2|X_1) = (\sqrt{1-0.0003} \cdot 9.9999) + (\sqrt{0.0003} \cdot 5) = 10.1344$$

$$q(X_3|X_2) = (\sqrt{1-0.0005} \cdot 1.0012) + (\sqrt{0.0005} \cdot 5) = \mathbf{10.24}$$

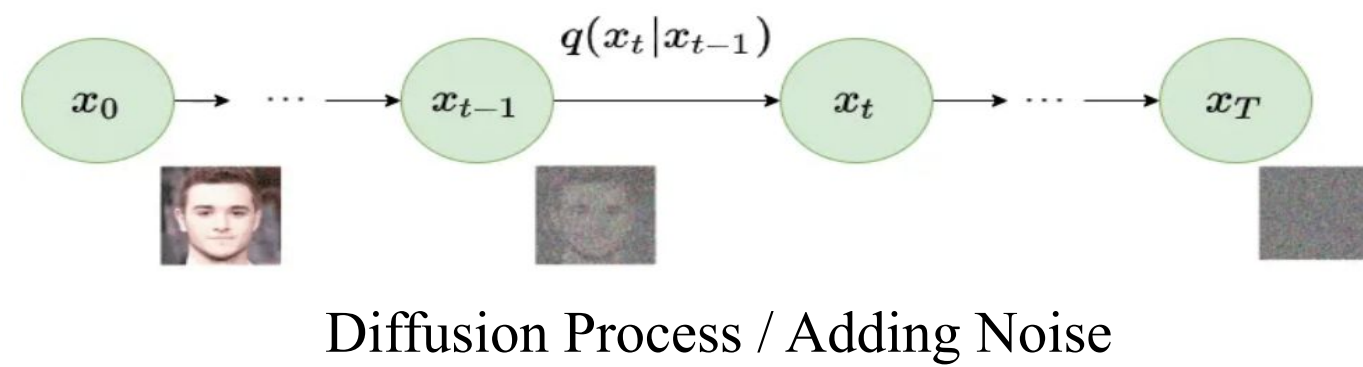
$$q(\mathbf{x}_t|\mathbf{x}_0) = \sqrt{\bar{\alpha}_t} \cdot \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \cdot \text{noise}$$

- $\alpha = 0.9999, \alpha = 0.9997, \alpha = 0.9995$
- $\alpha' = 0.9999, \alpha' = 0.9996, \alpha' = 0.9991$
- noise = 5
- input = 10

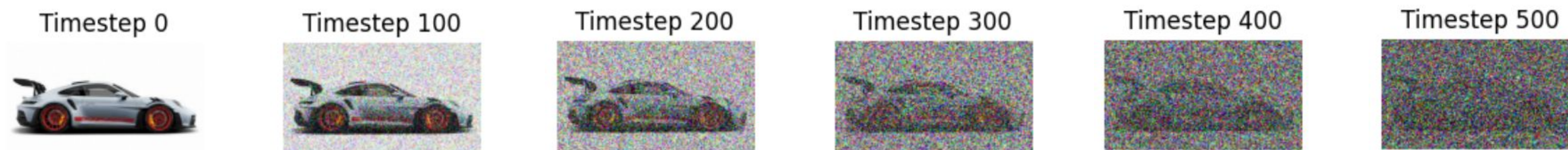
$$q(X_3|X_0) = \sqrt{0.9991} \cdot 10 + \sqrt{1-0.9991} \cdot 5 = \mathbf{10.14}$$



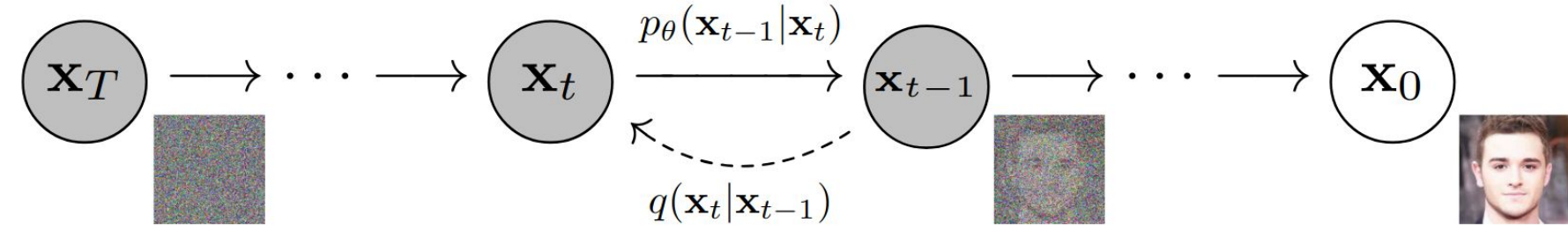
Diffusion Process / Adding Noise



$$q(\mathbf{x}_t|\mathbf{x}_0) = \sqrt{\bar{\alpha}_t} \cdot \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \cdot \text{noise}$$



Reverse Diffusion / Removing Noise



Reverse Diffusion Process / Removing Noise

$$p_{\theta}(\mathbf{x}_{0:T}) := p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t), \quad p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t))$$

This is what the model is trying to learn

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \cdot \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \cdot \mathbf{Z}$$

Where,

$$\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\sigma_t^2 = \beta_t$$

For more information please refer the paper



Overview

- Motivation behind diffusion models ✓
- Denoising Diffusion Probabilistic Models
 - Diffusion Process and Reverse Diffusion Process ✓
 - Training and Sampling a Diffusion Model
 - Results from a Diffusion Model that we trained
- Conditioned v/s Unconditioned Diffusion
- Policy Learning
- Toy Problem
- Diffusion Policy on Real Robots
 - Advantages
 - Network Architect
 - Evaluation & Results



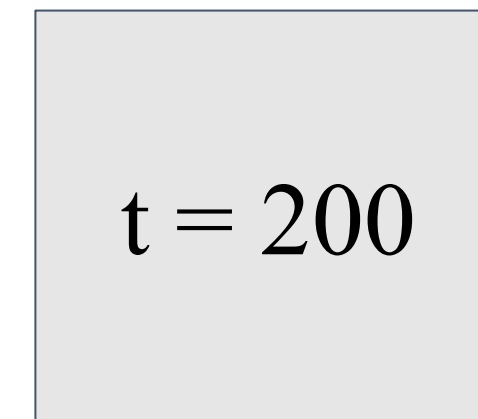
Training a Diffusion Model

Algorithm 1 Training

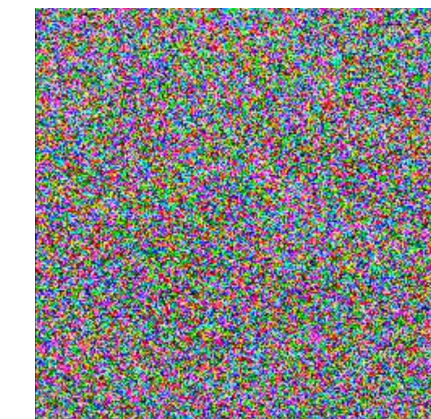
- 1: **repeat**
- 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
- 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
- 4: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 5: Take gradient descent step on $\nabla_{\theta} \|\epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2$
- 6: **until** converged



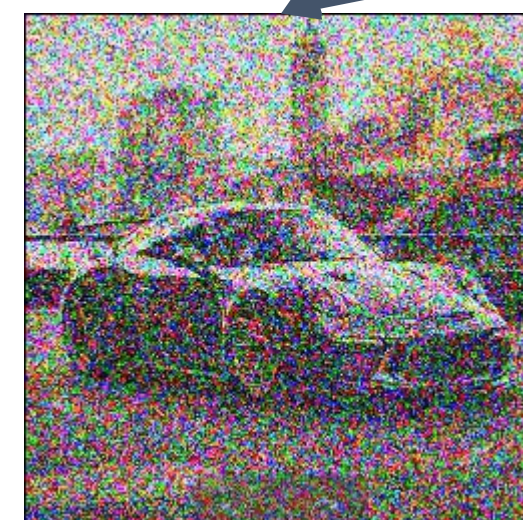
X_0



Random Time Step



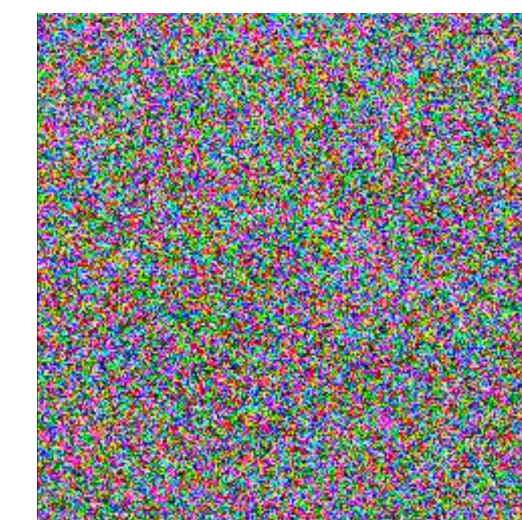
Random Noise (ϵ)



Noisy Image



UNet



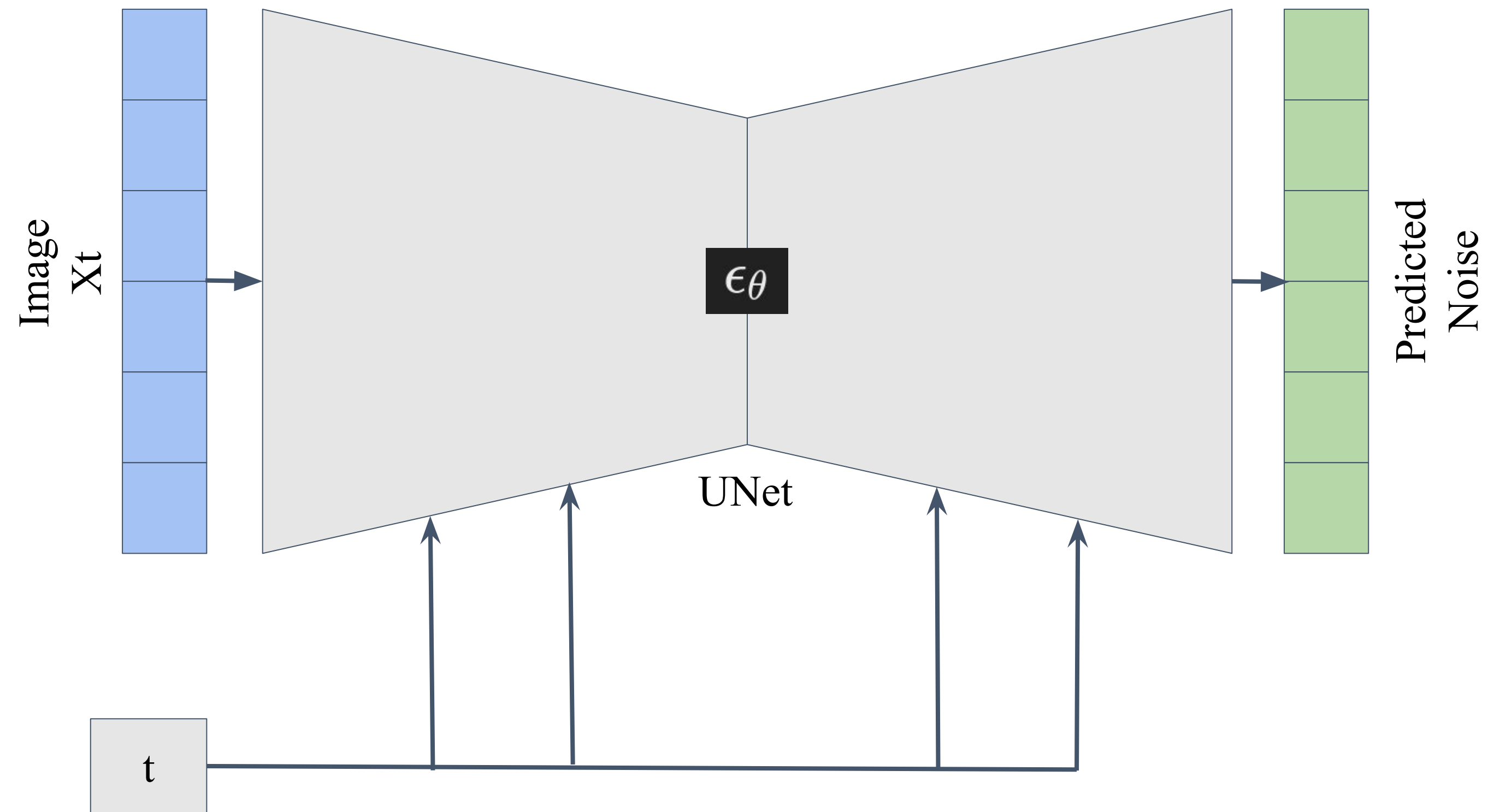
Predicted Noise

Time Step (t)



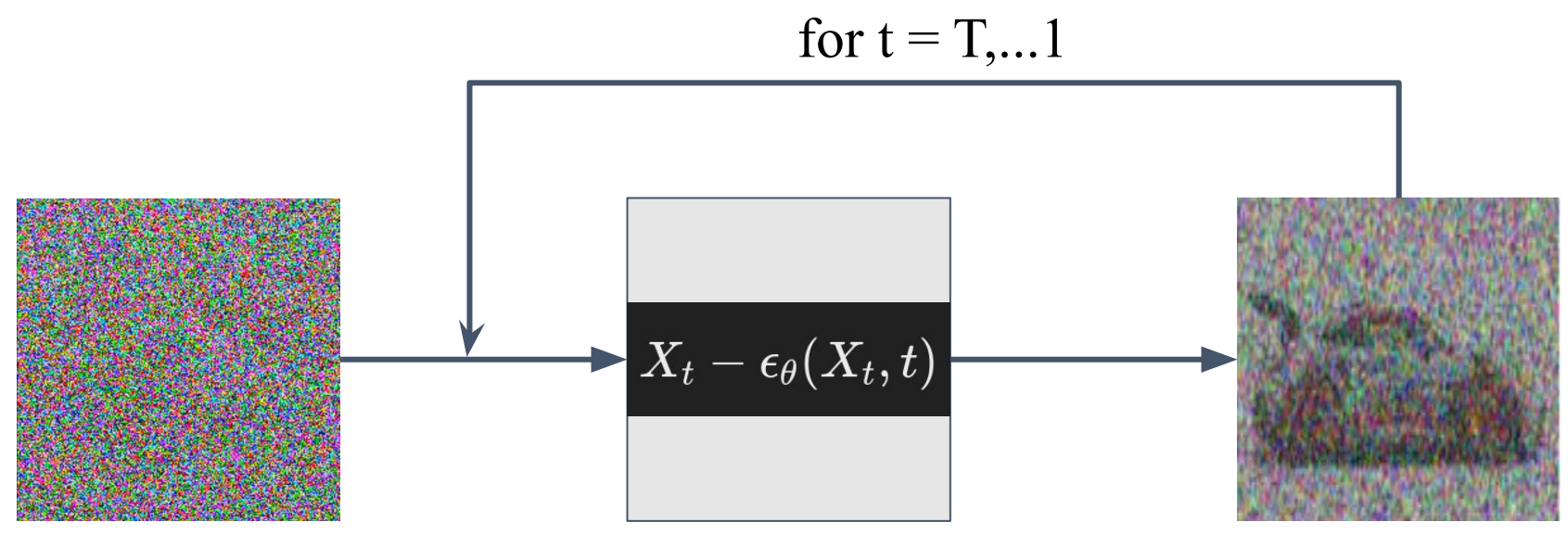
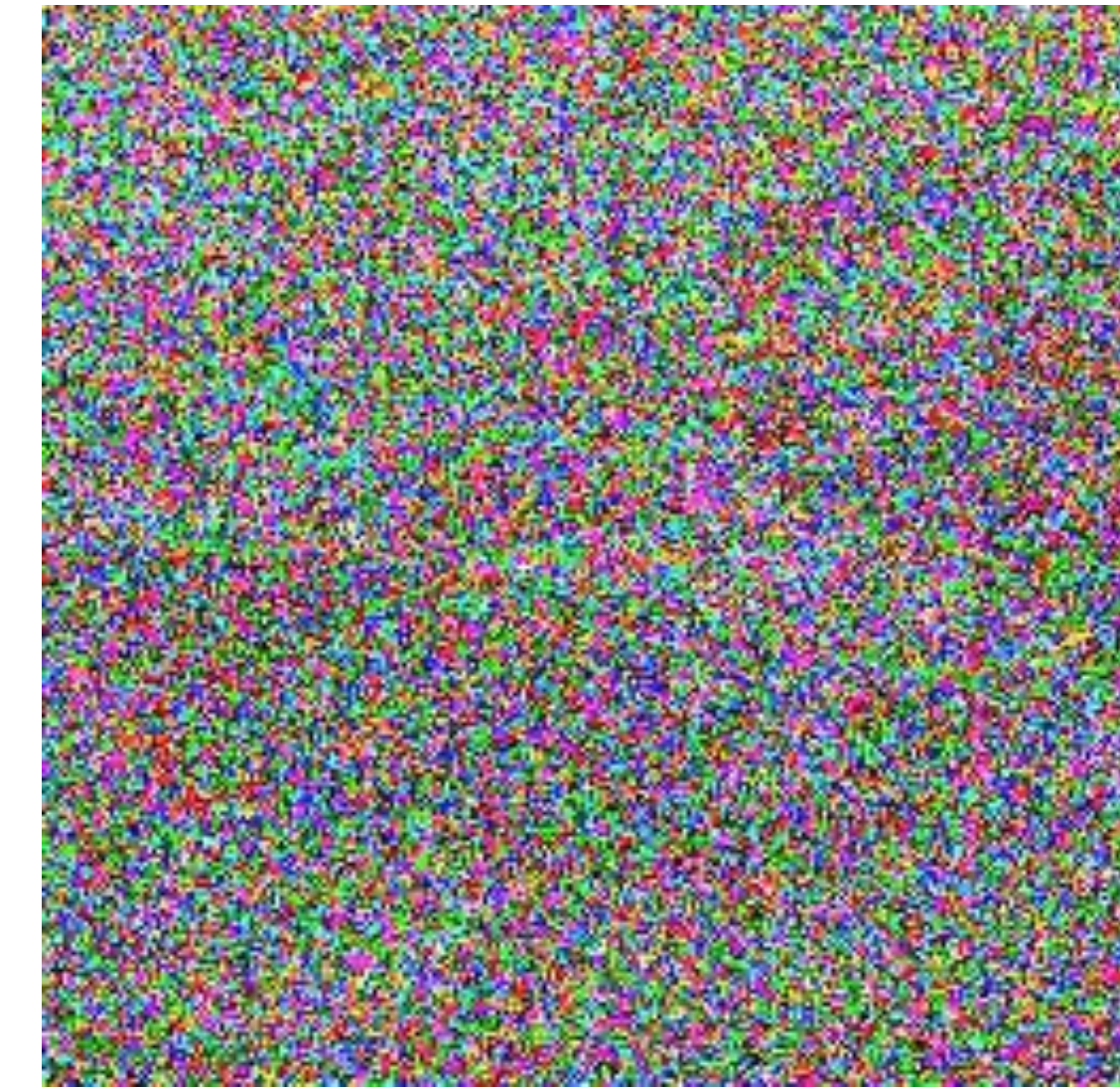


Training a Diffusion Model



Sampling from a Diffusion Model

```
Algorithm 2 Sampling  
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$   
2: for  $t = T, \dots, 1$  do  
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$   
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$   
5: end for  
6: return  $\mathbf{x}_0$ 
```



Overview

- Motivation behind diffusion models ✓
- Denoising Diffusion Probabilistic Models
 - Diffusion Process and Reverse Diffusion Process ✓
 - Training and Sampling a Diffusion Model ✓
 - Results from a Diffusion Model that we trained
- Conditioned v/s Unconditioned Diffusion
- Policy Learning
- Toy Problem
- Diffusion Policy
 - Advantages
 - Network Architect
 - Evaluation & Results





Denoising Diffusion Probabilistic Models



Stanford Cars Image Dataset



Diffusion Model

Hardware: RTX 3090 Ti
Train Time: 1.5 days
Epochs: 500



Denoising Diffusion Probabilistic Models



Diffused Cars!



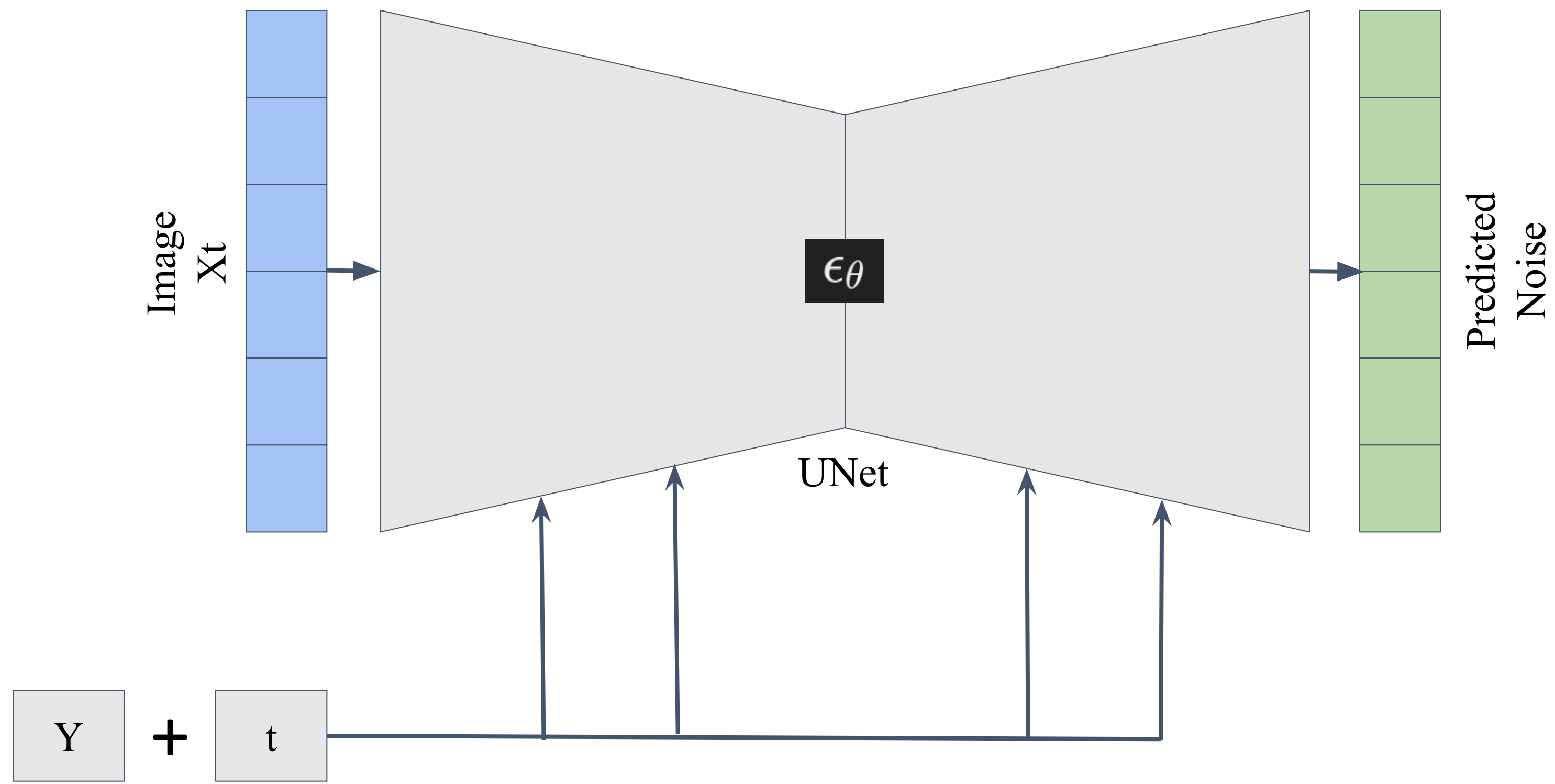
Overview

- Motivation behind diffusion models ✓
- Denoising Diffusion Probabilistic Models ✓
 - Diffusion Process and Reverse Diffusion Process ✓
 - Training and Sampling a Diffusion Model ✓
 - Results from a Diffusion Model that we trained ✓
- Conditioned v/s Unconditioned Diffusion
- Policy Learning
- Toy Problem
- Diffusion Policy on Real Robots
 - Advantages
 - Network Architect
 - Evaluation & Results

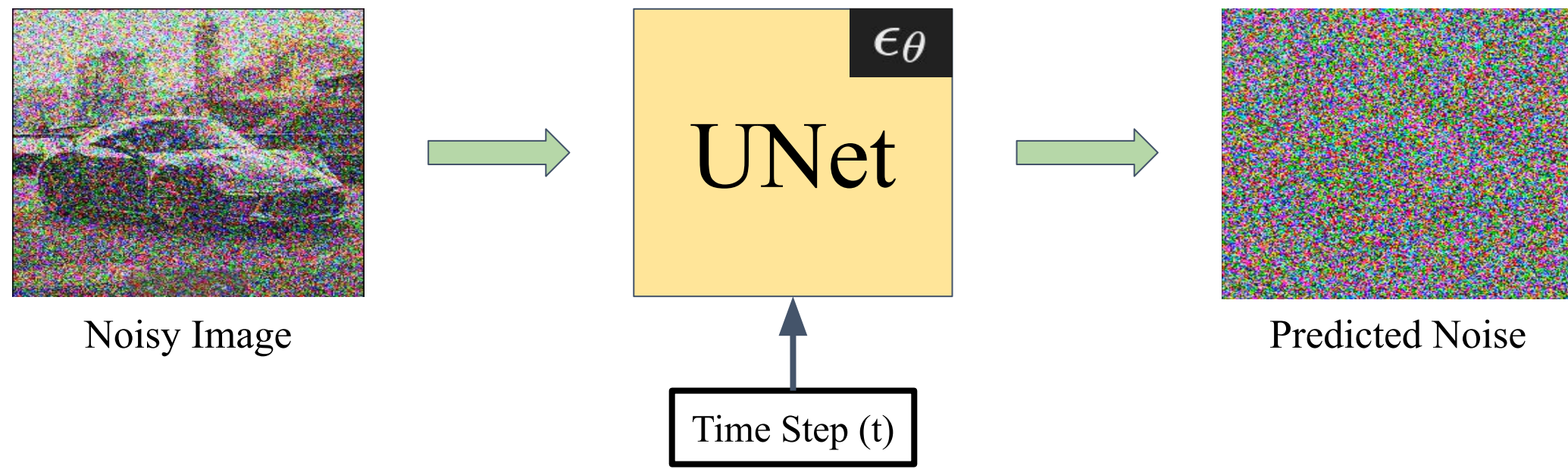




Training a Diffusion Model

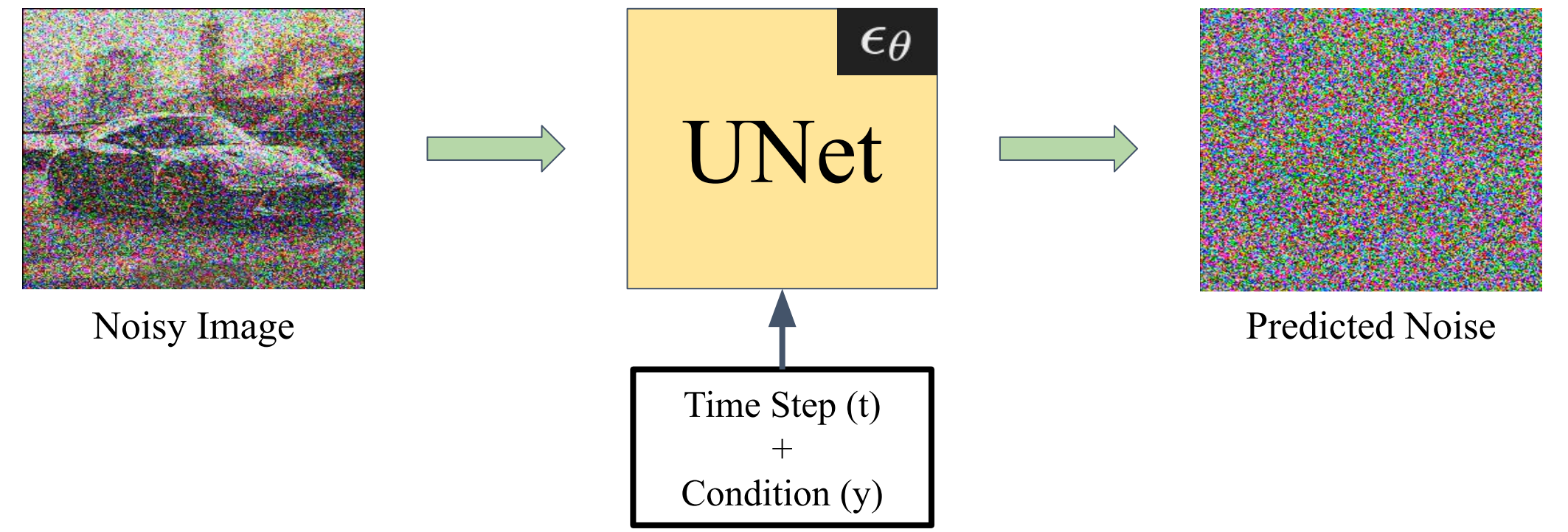


Unconditioned v/s Conditioned



$$\epsilon_{\theta}(\mathbf{x}_t, t)$$

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \cdot \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \cdot Z$$

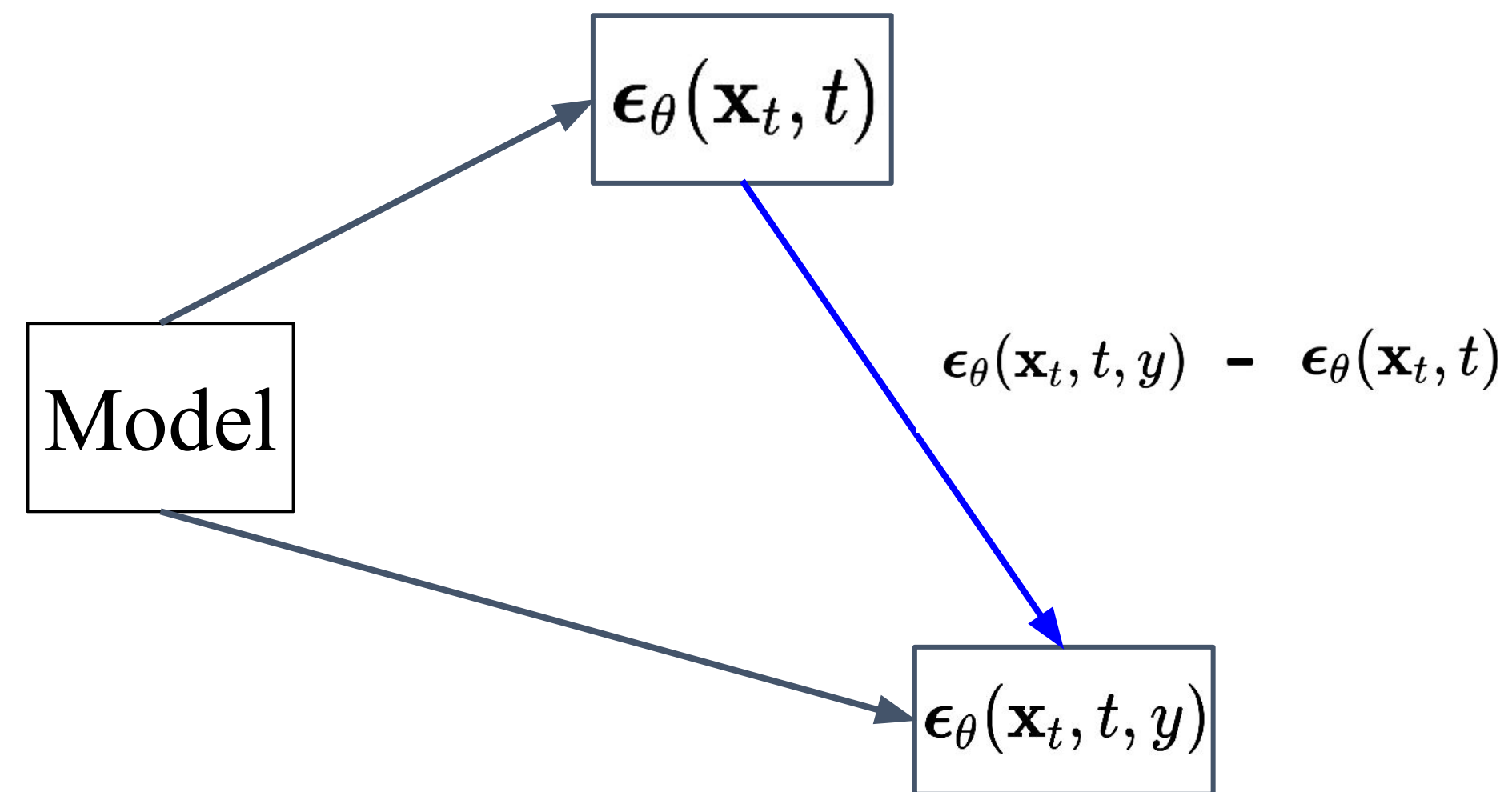


$$\epsilon_{\theta}(\mathbf{x}_t, t, y)$$

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \cdot \tilde{\epsilon}_{\theta}(\mathbf{x}_t, t, y) \right) + \sigma_t \cdot Z$$



Conditioned Diffusion



$$\tilde{\epsilon}_{\theta}(\mathbf{x}_t, t, y) = \epsilon_{\theta}(\mathbf{x}_t, t) + \lambda(\epsilon_{\theta}(\mathbf{x}_t, t, y) - \epsilon_{\theta}(\mathbf{x}_t, t))$$



Overview

- Motivation behind diffusion models ✓
- Denoising Diffusion Probabilistic Models ✓
 - Diffusion Process and Reverse Diffusion Process ✓
 - Training and Sampling a Diffusion Model ✓
 - Results from a Diffusion Model that we trained ✓
- Conditioned v/s Unconditioned Diffusion ✓
- Policy Learning
- Toy Problem
- Diffusion Policy on Real Robots
 - Advantages
 - Network Architect
 - Evaluation & Results





Policy Learning





Policy Learning

Policy: Given current observations, Policy tells a robot what to do next.





Policy Learning

Policy: Given current observations, Policy tells a robot what to do next.

$$A = \pi(O)$$



Policy Learning

Policy: Given current observations, Policy tells a robot what to do next.

$$A = \pi(O)$$



O



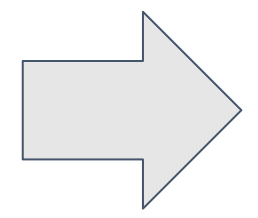
Policy Learning

Policy: Given current observations, Policy tells a robot what to do next.

$$A = \pi(O)$$



O

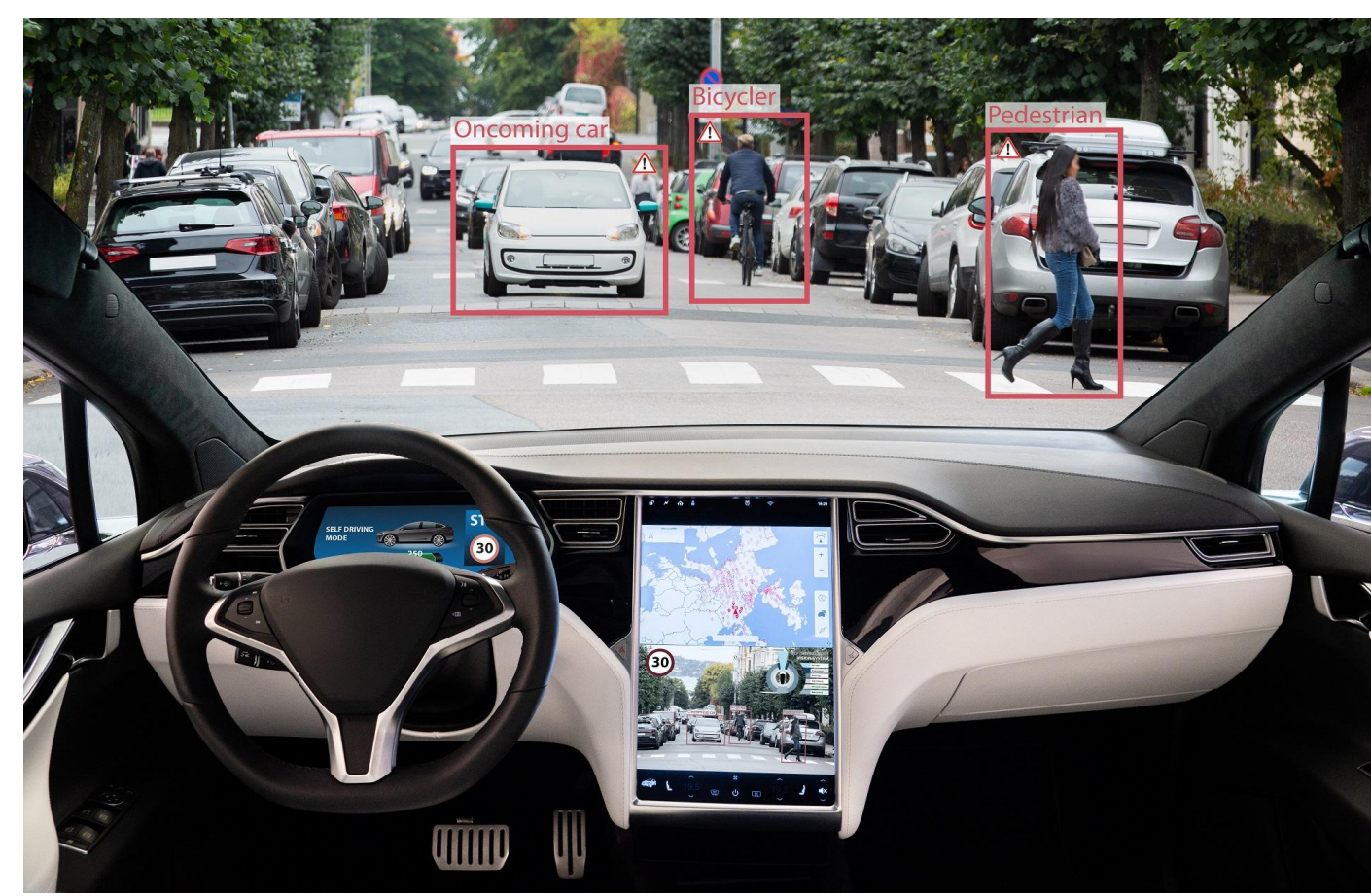


π

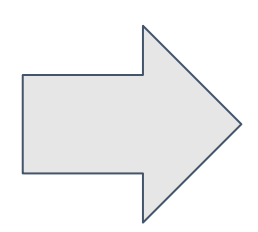
Policy Learning

Policy: Given current observations, Policy tells a robot what to do next.

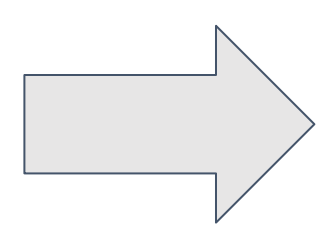
$$A = \pi(O)$$



O



π



“Turn Left”

A



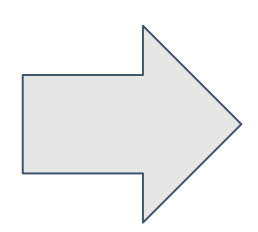
Policy Learning

Policy: Given current observations, Policy tells a robot what to do next.

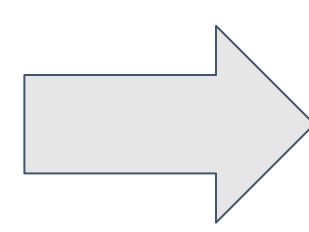
$$A = \pi(O)$$



O



π



“Turn Left”

A

Policy can be hardcoded or can be learned through data(**Policy Learning**)



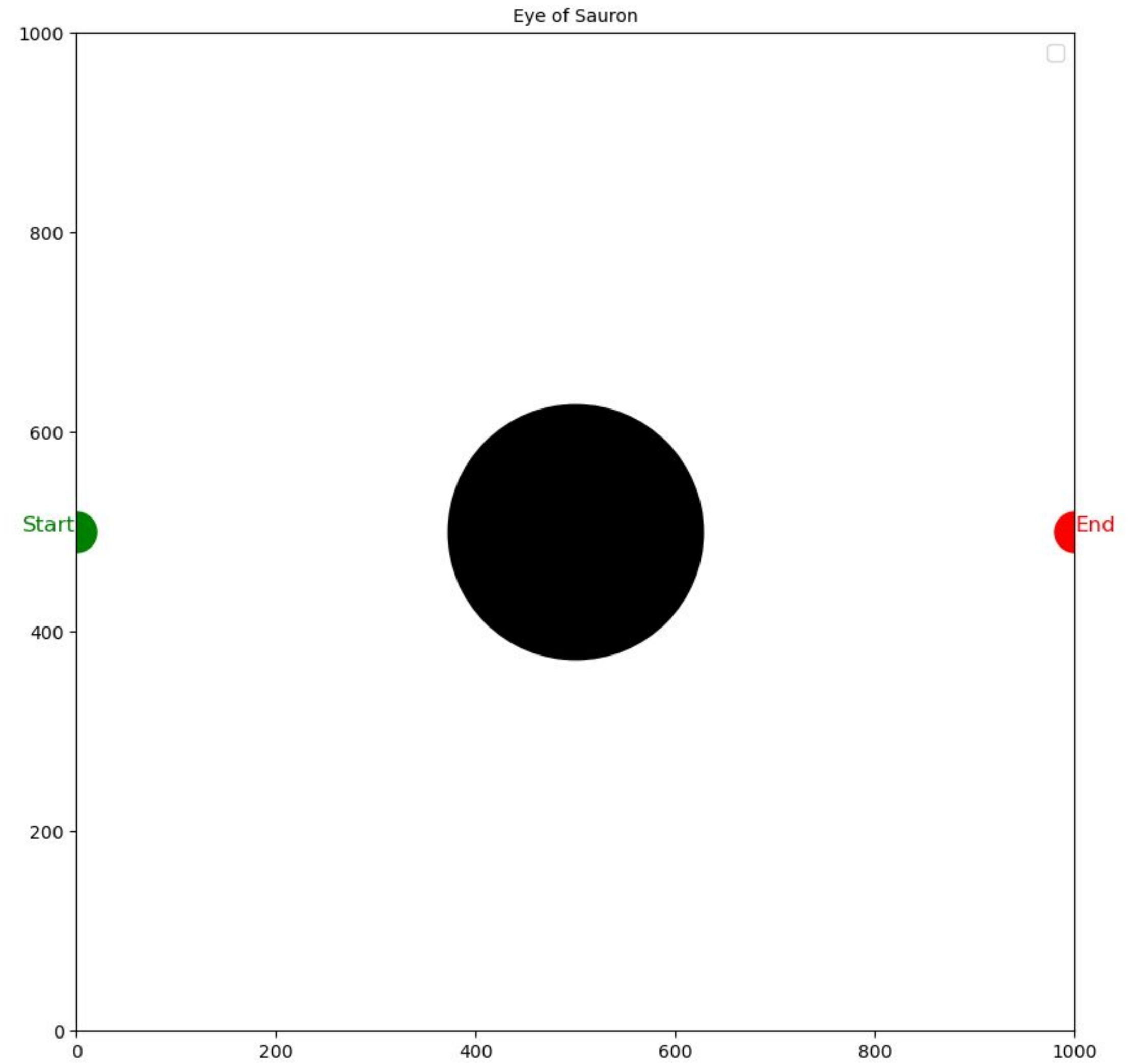


Toy Problem Definition





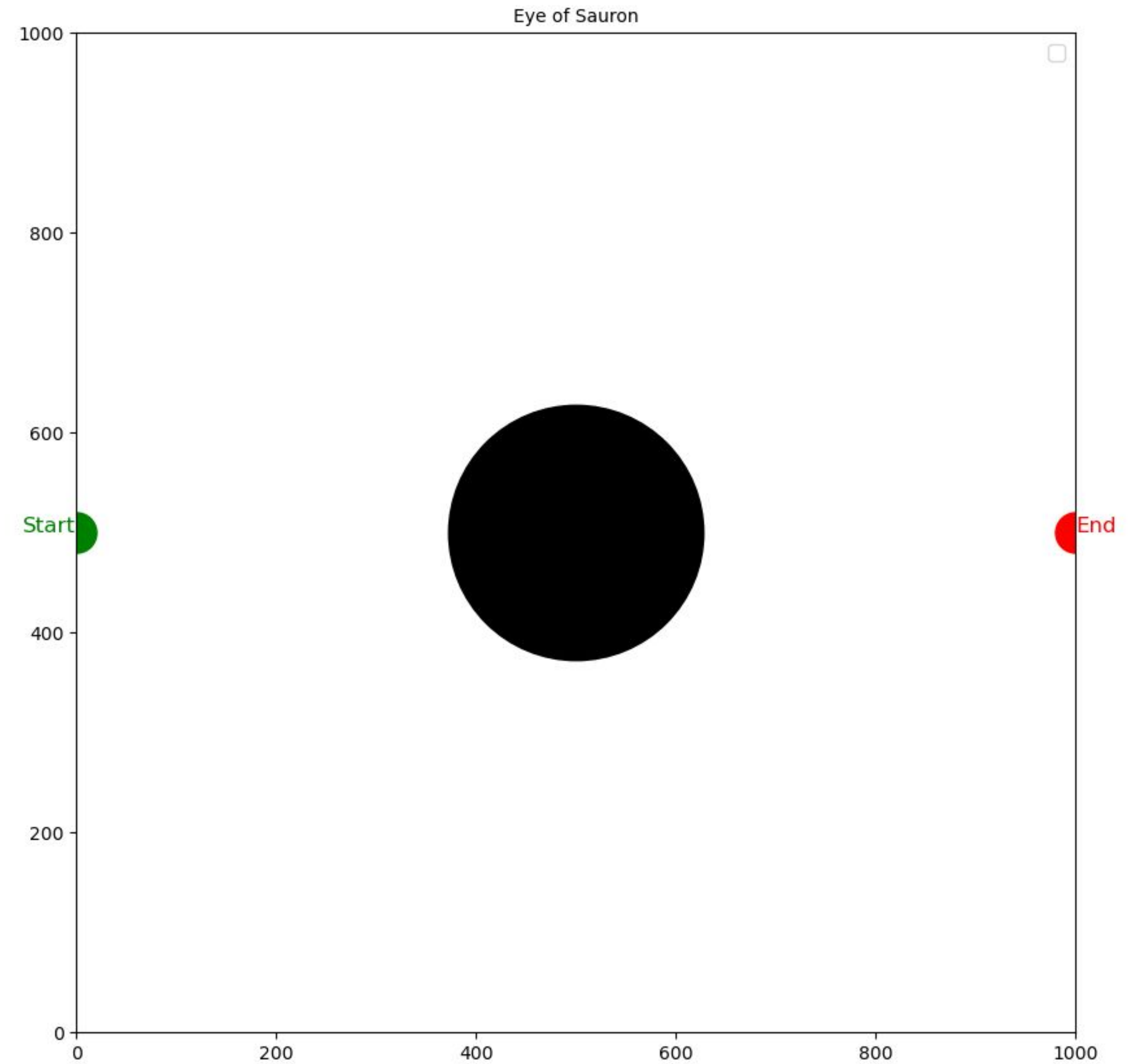
Toy Problem Definition





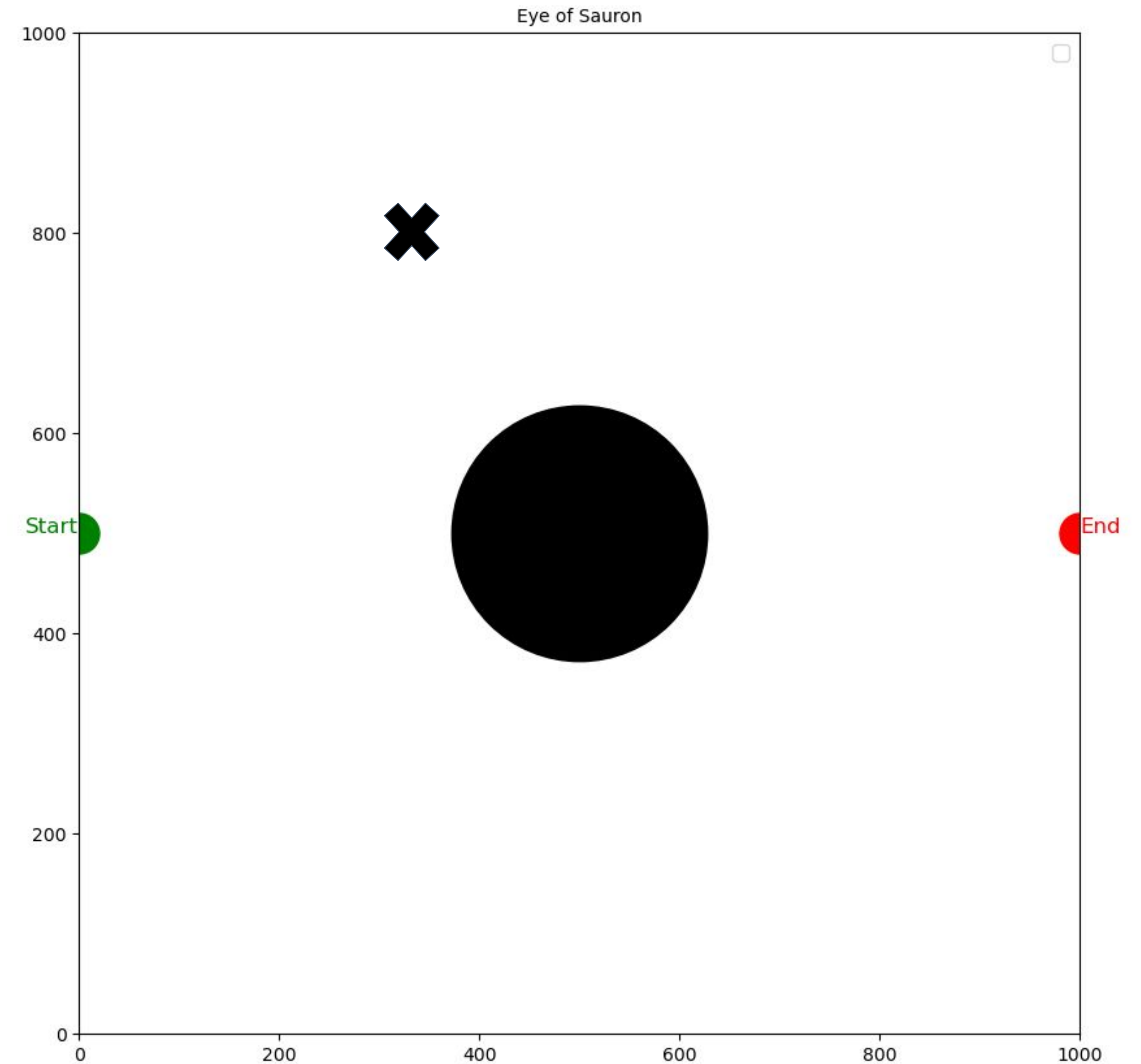
Toy Problem Definition

- Robot End Effector need to move from **Start** to the **End**



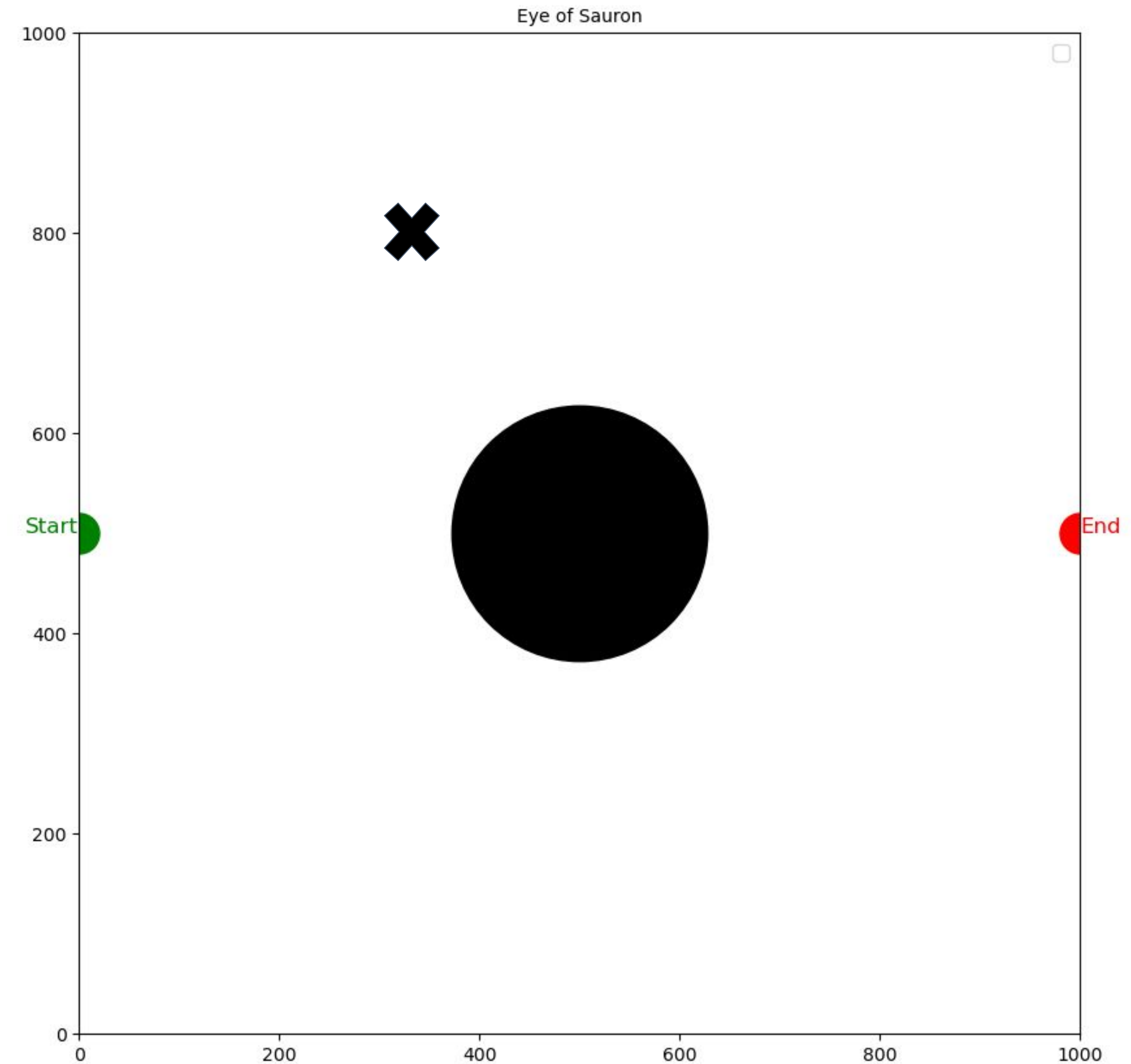
Toy Problem Definition

- Robot End Effector need to move from **Start** to the **End**.
- End Effector **✘** location(*State*) is defined by two variable $(x, y) : x, y \in [0, 1000]$
 - Example: EFF is at $(400, 800)$



Toy Problem Definition

- Robot End Effector need to move from **Start** to the **End**.
- End Effector **✘** location(*State*) is defined by two variable $(x, y) : x, y \in [0, 1000]$
 - Example: EFF is at $(400, 800)$
- Action is defined by two variables $(x, y) : x, y \in [0, 1000]$
 - Example: Go to $(405, 790)$





Toy Problem Data





Toy Problem Data

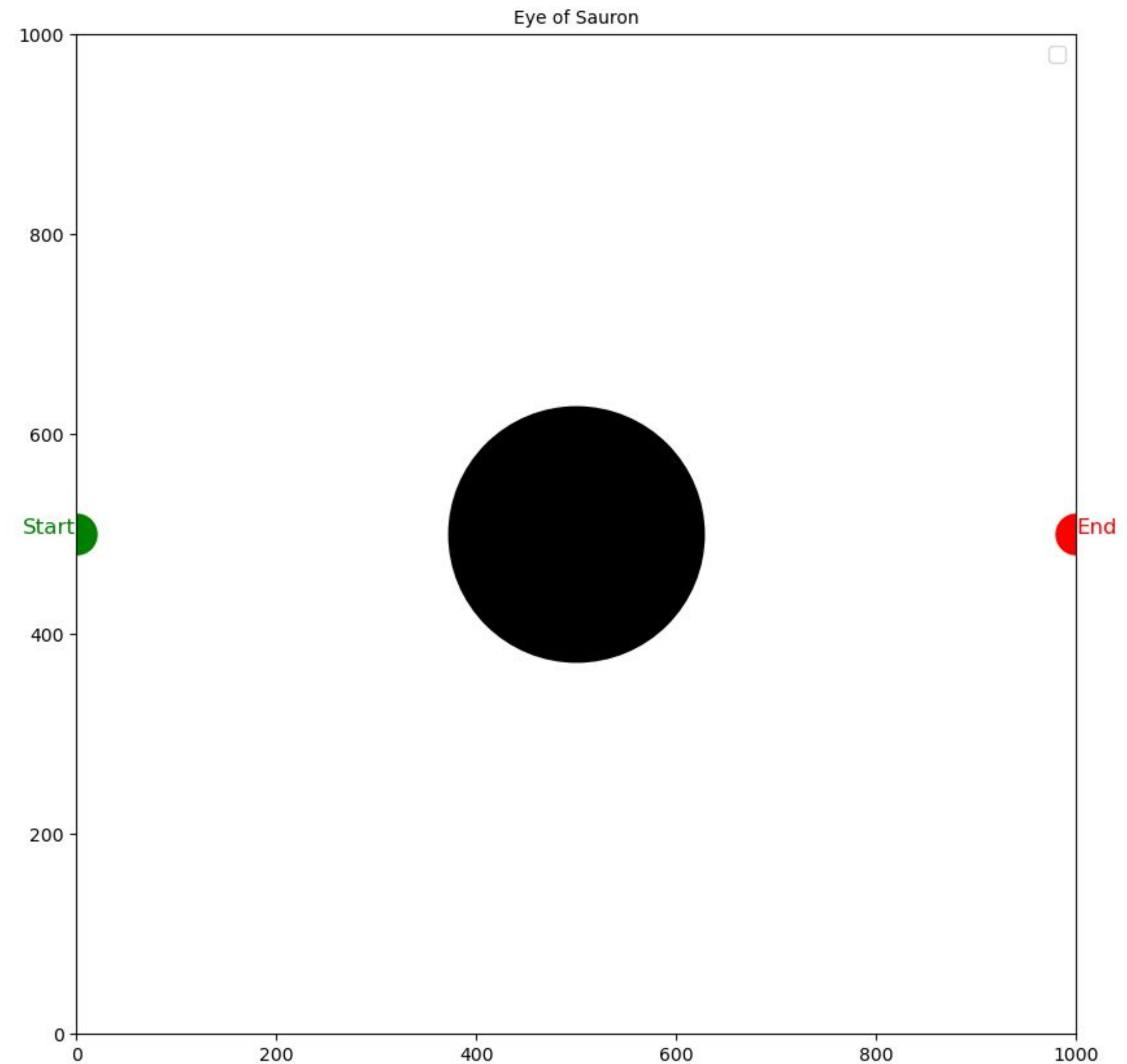
- We formulate our problem as a Supervised Learning problem.





Toy Problem Data

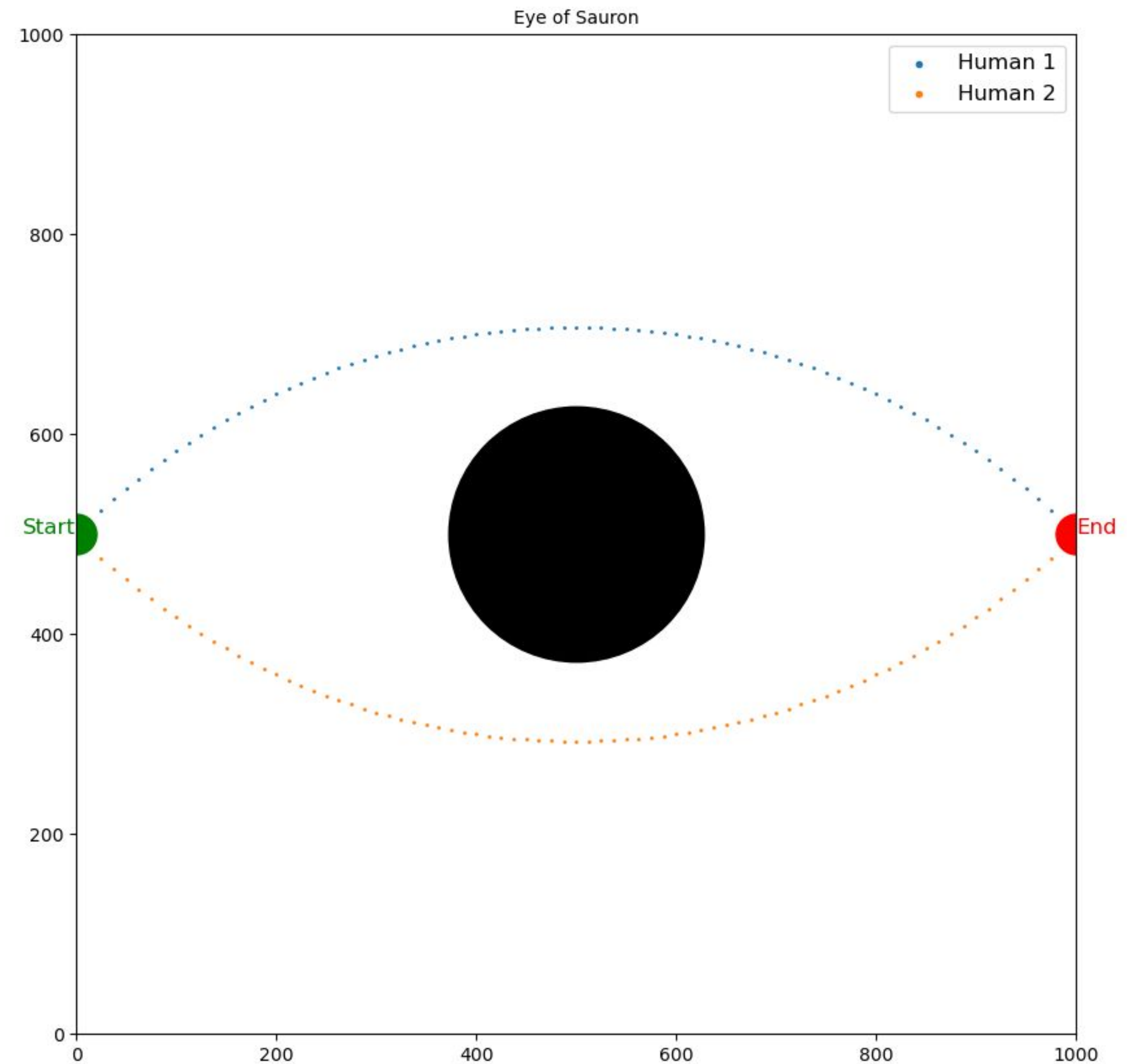
- We formulate our problem as a Supervised Learning problem.
- Let's say we collected data from human demonstration to complete the task





Toy Problem Data

- We formulate our problem as a Supervised Learning problem.
- Let's say we collected data from human demonstration to complete the task***



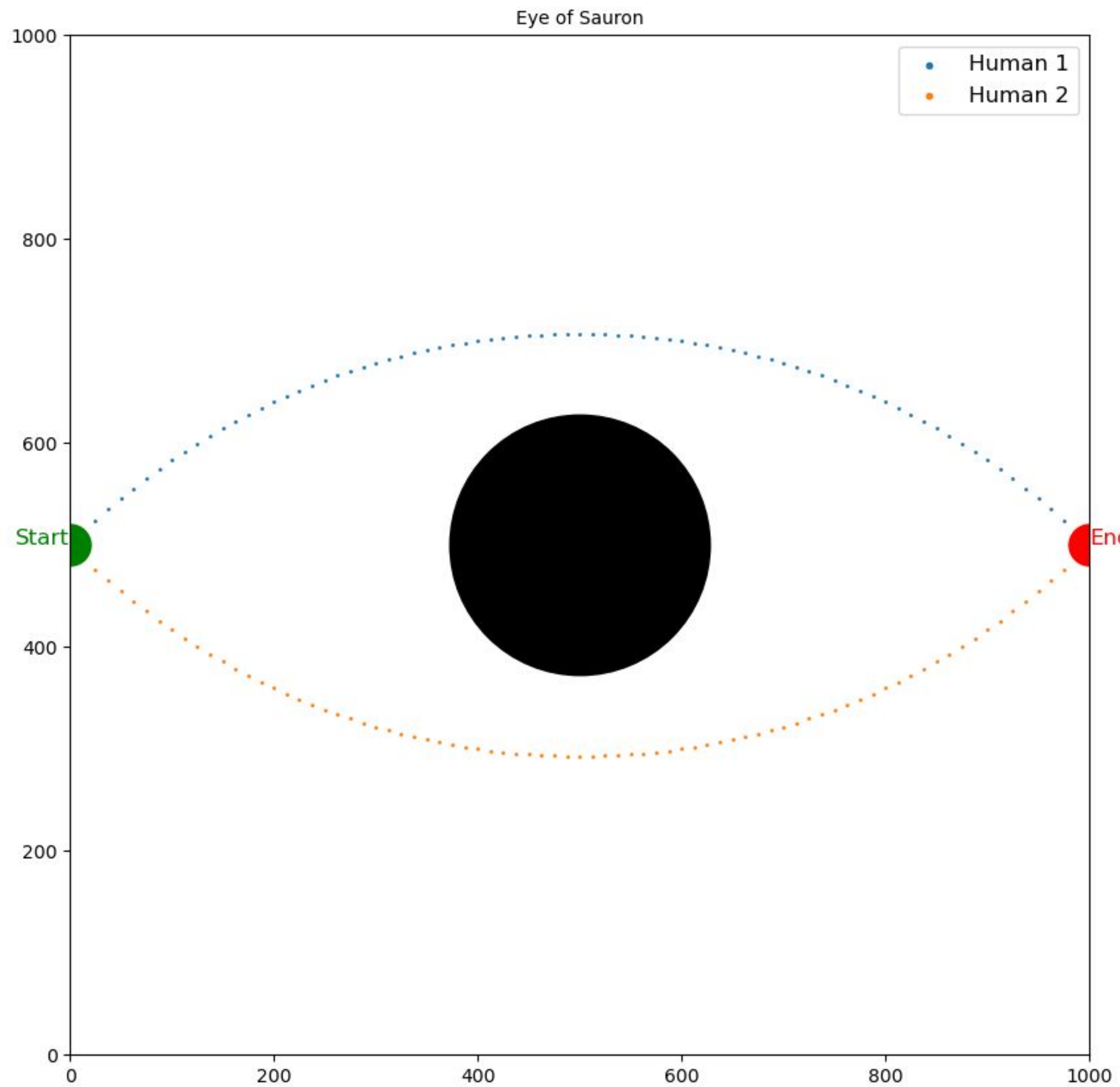
*** toy data was generated using two equations of circle centered at (500,0) and (500,1000).





Toy Problem Data

- We formulate our problem as a Supervised Learning problem.
- Let's say we collected data from human demonstration to complete the task***
- Every dot in the plot represents action taken by robot's end effector



*** toy data was generated using two equations of circle centered at (500,0) and (500,1000).





Toy Problem Data

Now we have “State” to “Action” mapping to train our Model.





Toy Problem Data

Now we have “State” to “Action” mapping to train our Model.

State/
Observation

X	0	12	25	38	50	62	75	88	100	112	125	138	150	162	175
Y	500	512	524	535	545	556	565	574	583	591	599	607	614	621	628



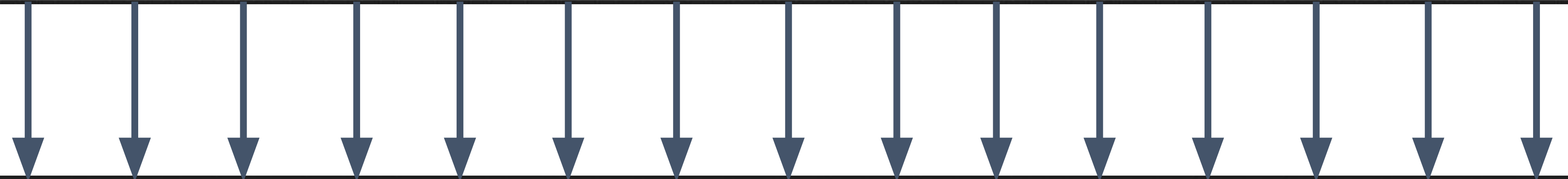


Toy Problem Data

Now we have “State” to “Action” mapping to train our Model.

State/
Observation

X	0	12	25	38	50	62	75	88	100	112	125	138	150	162	175
Y	500	512	524	535	545	556	565	574	583	591	599	607	614	621	628



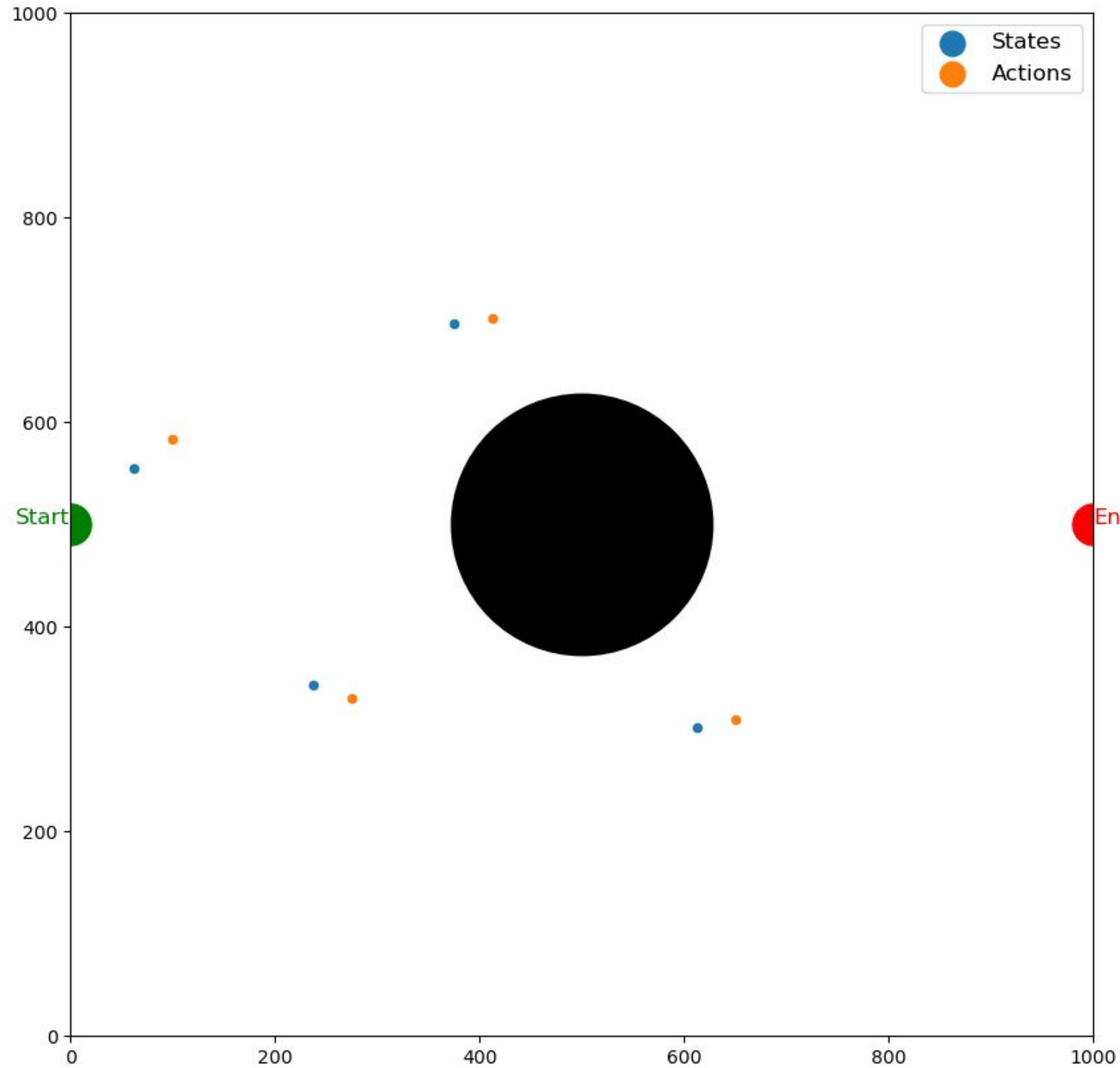
Action

X	38	50	62	75	88	100	112	125	138	150	162	175	188	200	212
Y	535	545	556	565	574	583	591	599	607	614	621	628	634	640	646



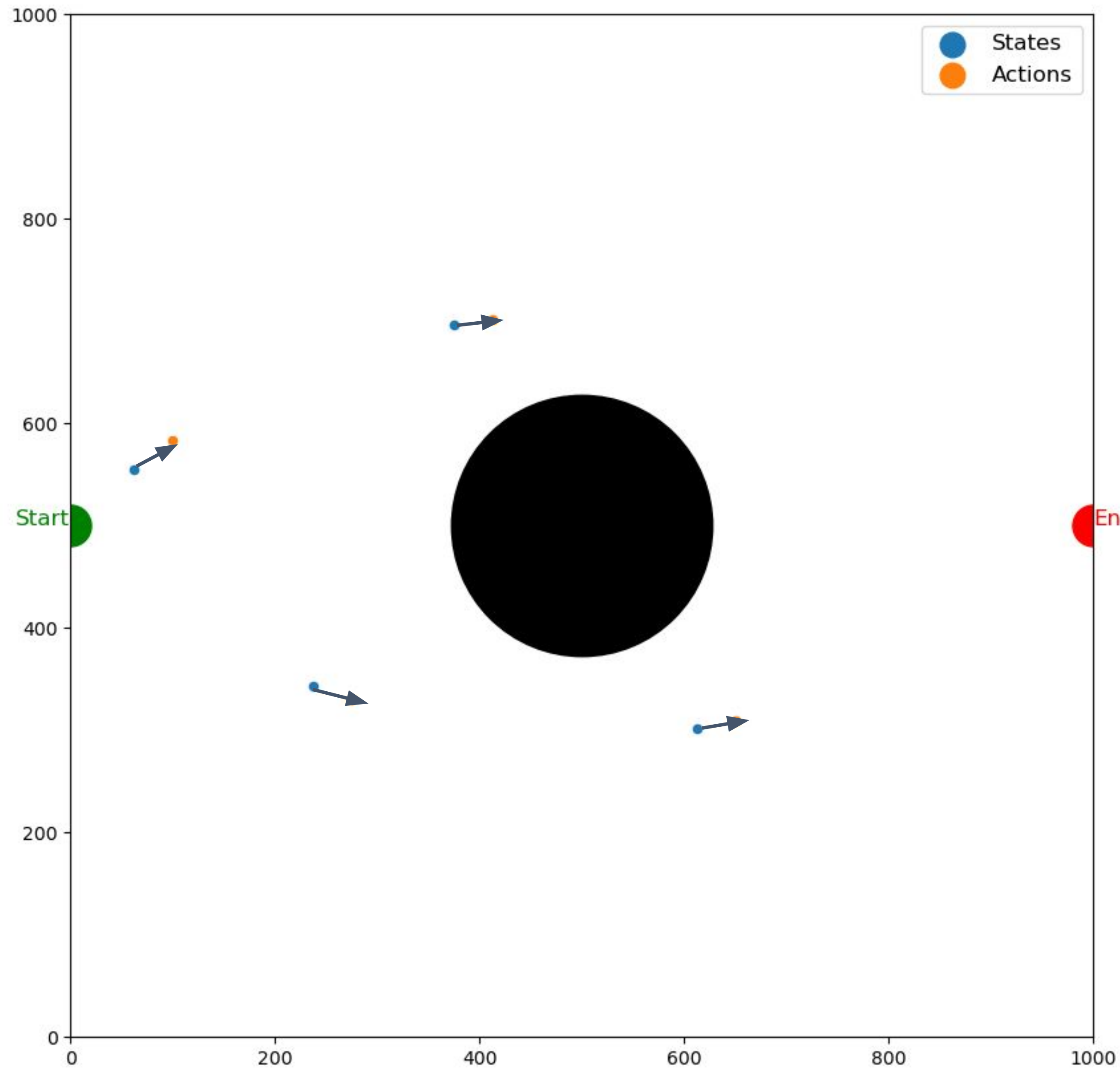


Toy Problem Data





Toy Problem Data





Some Nuances before Building a Model





Some Nuances before Building a Model

- It takes time to infer from the model, robot will be idle for a this time?





Some Nuances before Building a Model

- It takes time to infer from the model, robot will be idle for a this time?
- **Move the robot while model is predicting in parallel.**





Some Nuances before Building a Model

- It takes time to infer from the model, robot will be idle for a this time?
- **Move the robot while model is predicting in parallel.**
- Only having information about current location is enough??



Some Nuances before Building a Model

- It takes time to infer from the model, robot will be idle for a this time?
- **Move the robot while model is predicting in parallel.**
- Only having information about current location is enough??
- **Provide history of states**





Some Nuances before Building a Model

- It takes time to infer from the model, robot will be idle for a this time?
- **Move the robot while model is predicting in parallel.**
- Only having information about current location is enough??
- **Provide history of states**
- Only predicting the next Action is not efficient, motion will be Zig-Zaggy





Some Nuances before Building a Model

- It takes time to infer from the model, robot will be idle for a this time?
- **Move the robot while model is predicting in parallel.**
- Only having information about current location is enough??
- **Provide history of states**
- Only predicting the next Action is not efficient, motion will be Zig-Zaggy
- **Predict sequence of actions**





Some Nuances before Building a Model

- For Toy Example:
 - Observation Horizon: 5
 - Prediction Horizon: 16
 - Action Horizon: 8 (used when policy is rolled out)

	t																		
Observation:	█	█	█	█	█														
Action Pred:				█	█	█	█	█	█	█	█	█	█	█	█	█	█	█	█
Action Used:					█	█	█	█	█	█	█								
time:	-4	-3	-2	-1	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14





Toy Model Denoising

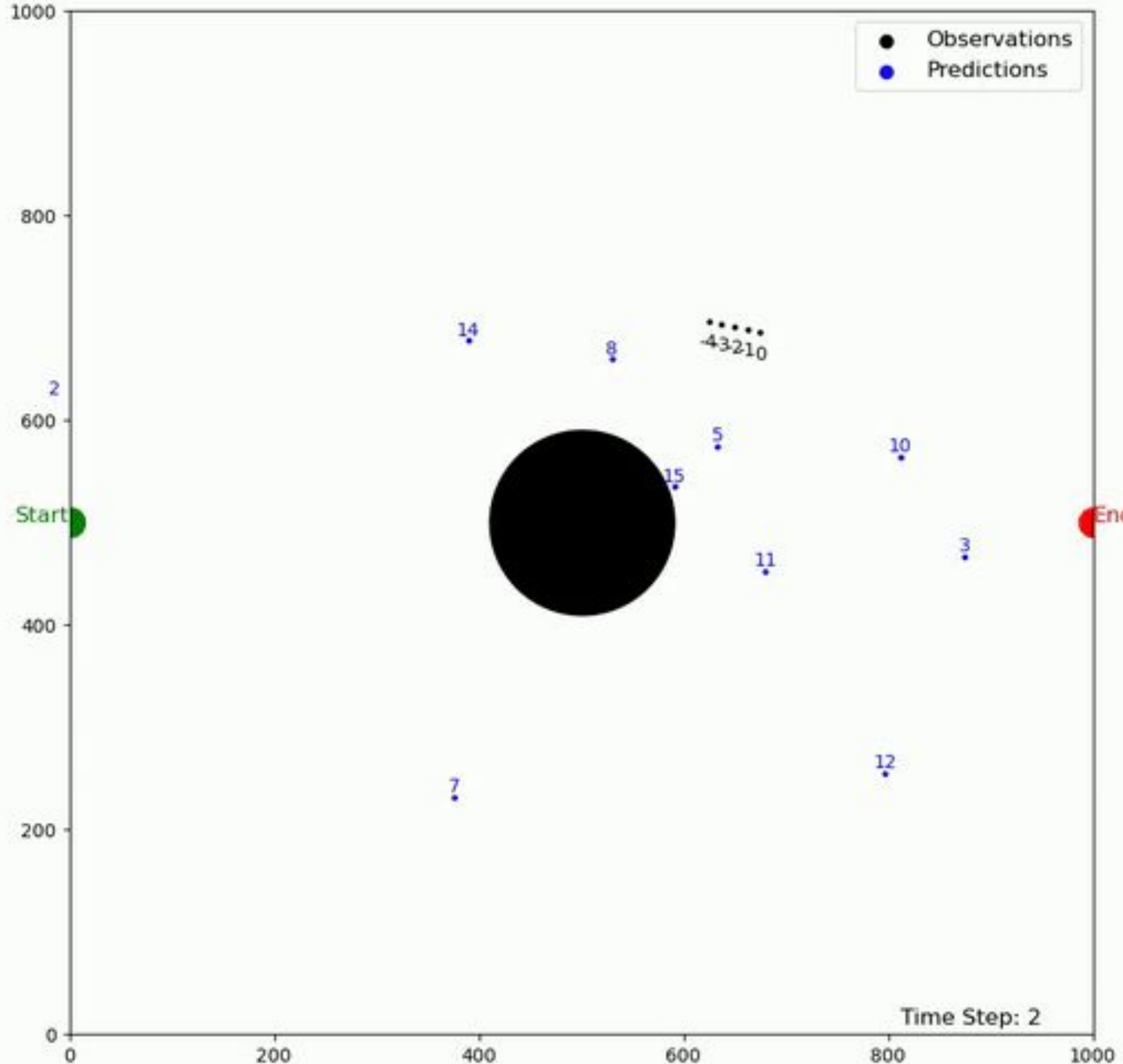
With the setup discuss we train a Diffusion Model to predict next 16 actions conditioned upon the last 5 states.

The output from the model is visualized in next slide.





Toy Model Denoising



Overview

- Motivation behind diffusion models ✓
- Denoising Diffusion Probabilistic Models ✓
 - Diffusion Process and Reverse Diffusion Process ✓
 - Training and Sampling a Diffusion Model ✓
 - Results from a Diffusion Model that we trained ✓
- Conditioned v/s Unconditioned Diffusion ✓
- Policy Learning ✓
- Toy Problem ✓
- Diffusion Policy on Real Robots
 - Advantages
 - Network Architect
 - Evaluation & Results





Transfer to Real Robot for Manipulation

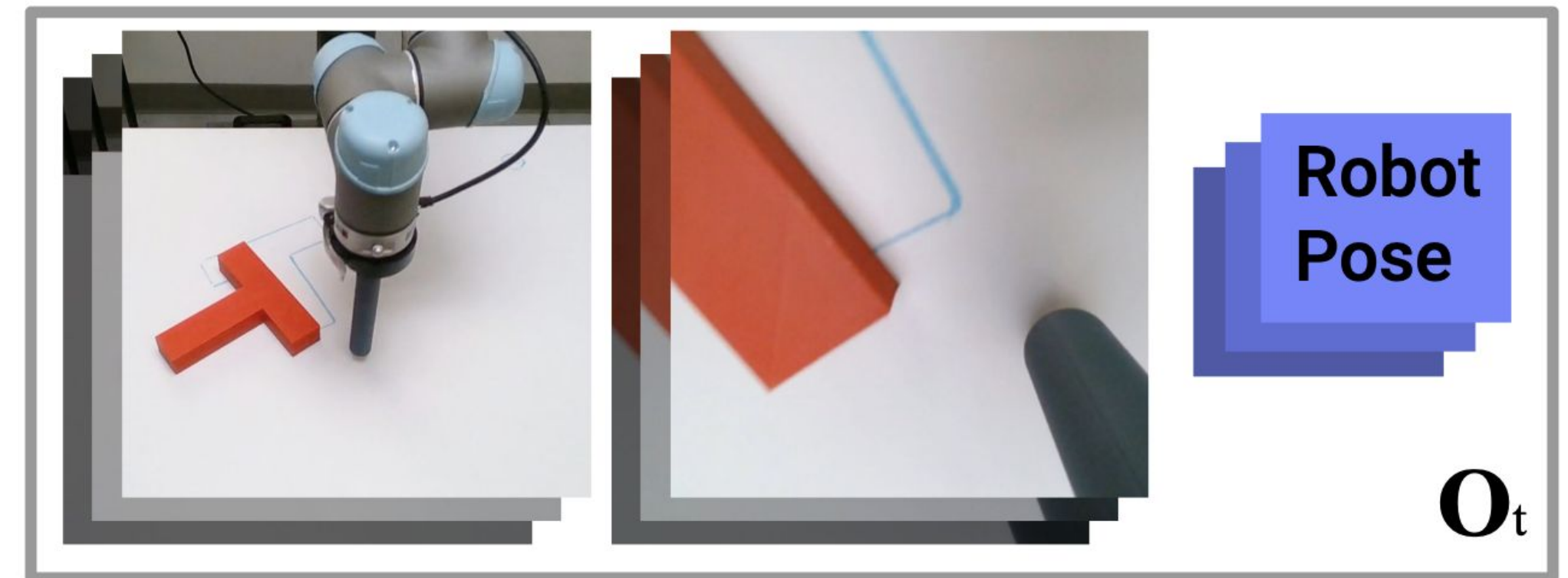


Transfer to Real Robot for Manipulation

- Observation could include
 - Robot EEF location & orientation
 - Robot Joint Configuration
 - Image Input
 - Other information that you want your model to have....

- Action could be:
 - Robot EEF location & orientation
 - Robot Joint Configuration
 - Joint Velocities, etc.

Image Observation Sequence





Diffusion Policy Paper

How does the Maths transfer from Image Space?



How does Maths Transfer from Image Space?

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



How does Maths Transfer from Image Space?

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

Conditional Image Generation $\mathbf{x}^{k-1} = \alpha(\mathbf{x}^k - \gamma \varepsilon_{\theta}(\mathbf{x}^k, \mathbf{y}, k) + \mathcal{N}(0, \sigma^2 I))$



How does Maths Transfer from Image Space?

Two Modifications:

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

Conditional Image Generation $\mathbf{x}^{k-1} = \alpha(\mathbf{x}^k - \gamma \varepsilon_{\theta}(\mathbf{x}^k, \mathbf{y}, k) + \mathcal{N}(0, \sigma^2 I))$



How does Maths Transfer from Image Space?

Two Modifications:

1. Image(x) -> Action(A)

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

Conditional Image Generation $\mathbf{x}^{k-1} = \alpha(\mathbf{x}^k - \gamma \varepsilon_{\theta}(\mathbf{x}^k, \mathbf{y}, k) + \mathcal{N}(0, \sigma^2 I))$



How does Maths Transfer from Image Space?

Two Modifications:

1. Image(x) -> Action(A)
2. Denoising conditioned on observation(O)

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

Conditional Image Generation $\mathbf{x}^{k-1} = \alpha(\mathbf{x}^k - \gamma \varepsilon_{\theta}(\mathbf{x}^k, \mathbf{y}, k) + \mathcal{N}(0, \sigma^2 I))$



How does Maths Transfer from Image Space?

Two Modifications:

1. Image(x) -> Action(A)
2. Denoising conditioned on observation(O)

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

Conditional Image Generation $\mathbf{x}^{k-1} = \alpha(\mathbf{x}^k - \gamma \varepsilon_{\theta}(\mathbf{x}^k, \mathbf{y}, k) + \mathcal{N}(0, \sigma^2 I))$

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



How does Maths Transfer from Image Space?

Two Modifications:

1. Image(x) -> Action(A)
2. Denoising conditioned on observation(O)

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

Conditional Image Generation $\mathbf{x}^{k-1} = \alpha(\mathbf{x}^k - \gamma \varepsilon_{\theta}(\mathbf{x}^k, \mathbf{y}, k) + \mathcal{N}(0, \sigma^2 I))$

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





Diffusion Policy Paper

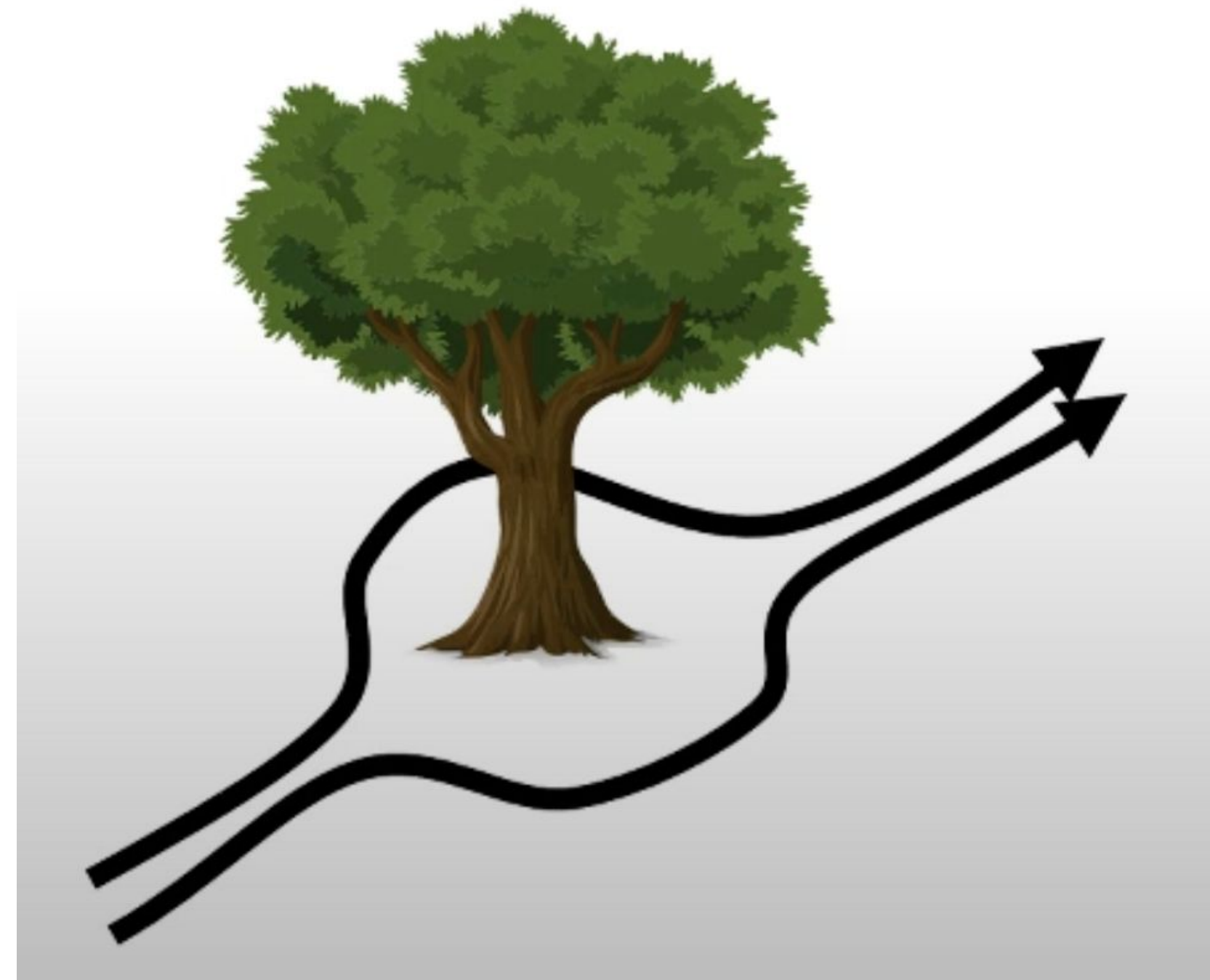
What advantages does using diffusion for learning policy have?





1. MultiModal Action Distribution

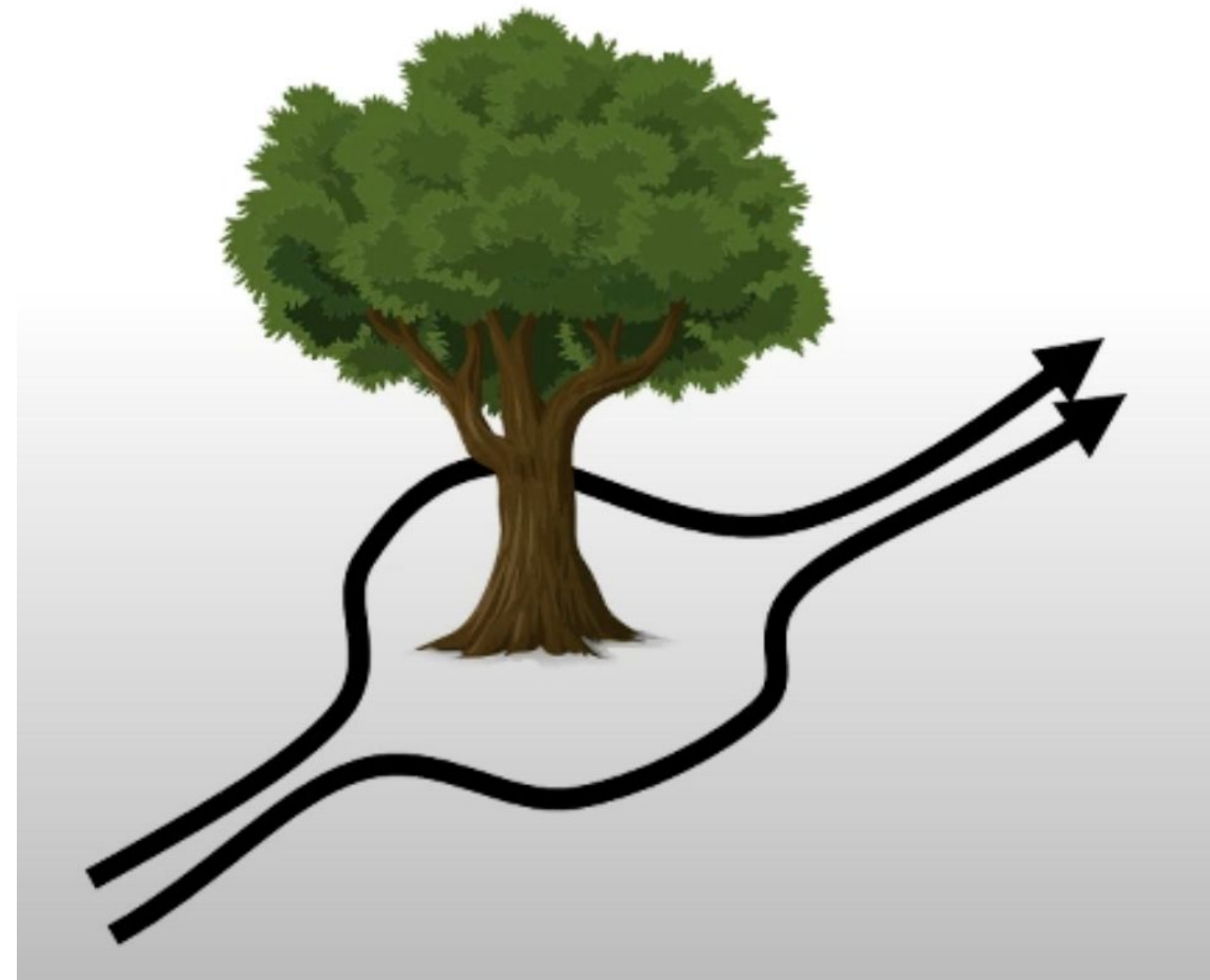
If half of the data say move right and the other half says move left:





1. MultiModal Action Distribution

If half of the data say move right and the other half says move left:
Simple MLPs -> Move Center



1. MultiModal Action Distribution

If half of the data say move right and the other half says move left:

Simple MLPs -> Move Center

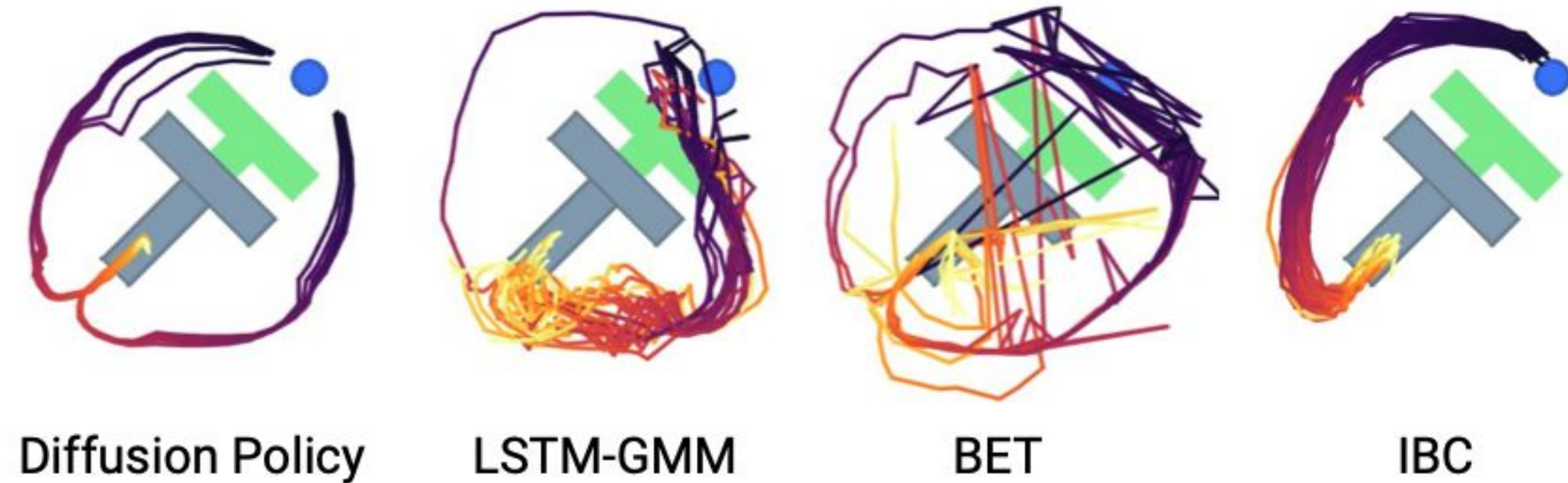


Diffusion Policy



1. MultiModal Action Distribution

If half of the data say move right and the other half says move left:
Simple MLPs -> Move Center





1. MultiModal Action Distribution

If half of the data say move right and the other half says move left:
Simple MLPs -> Move Center

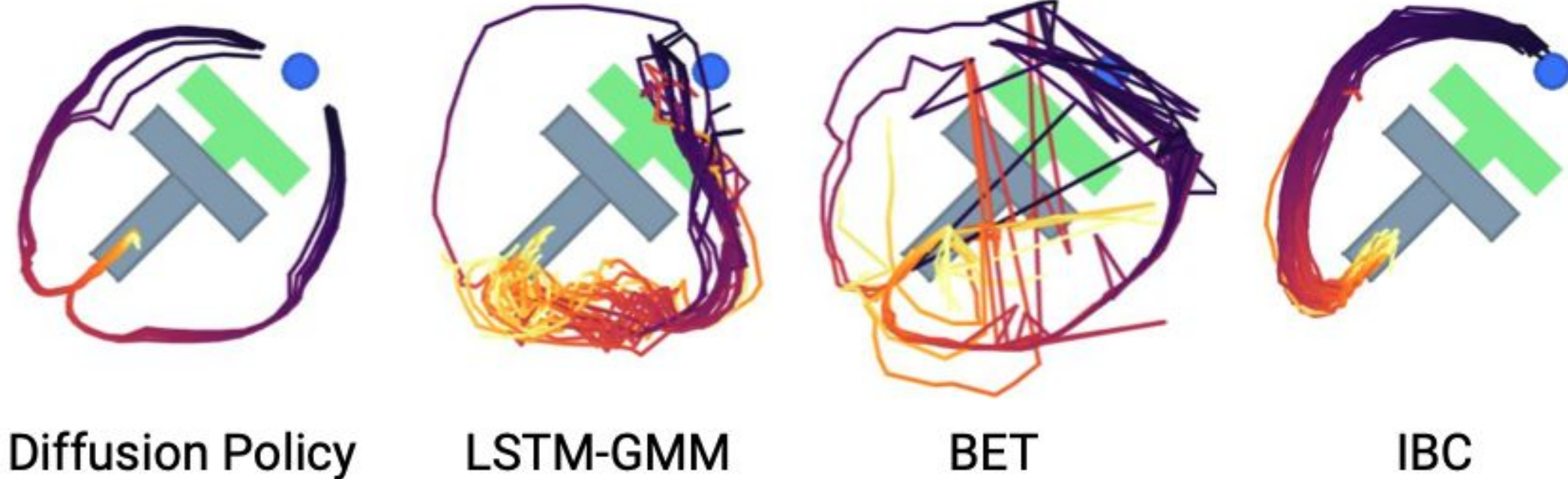
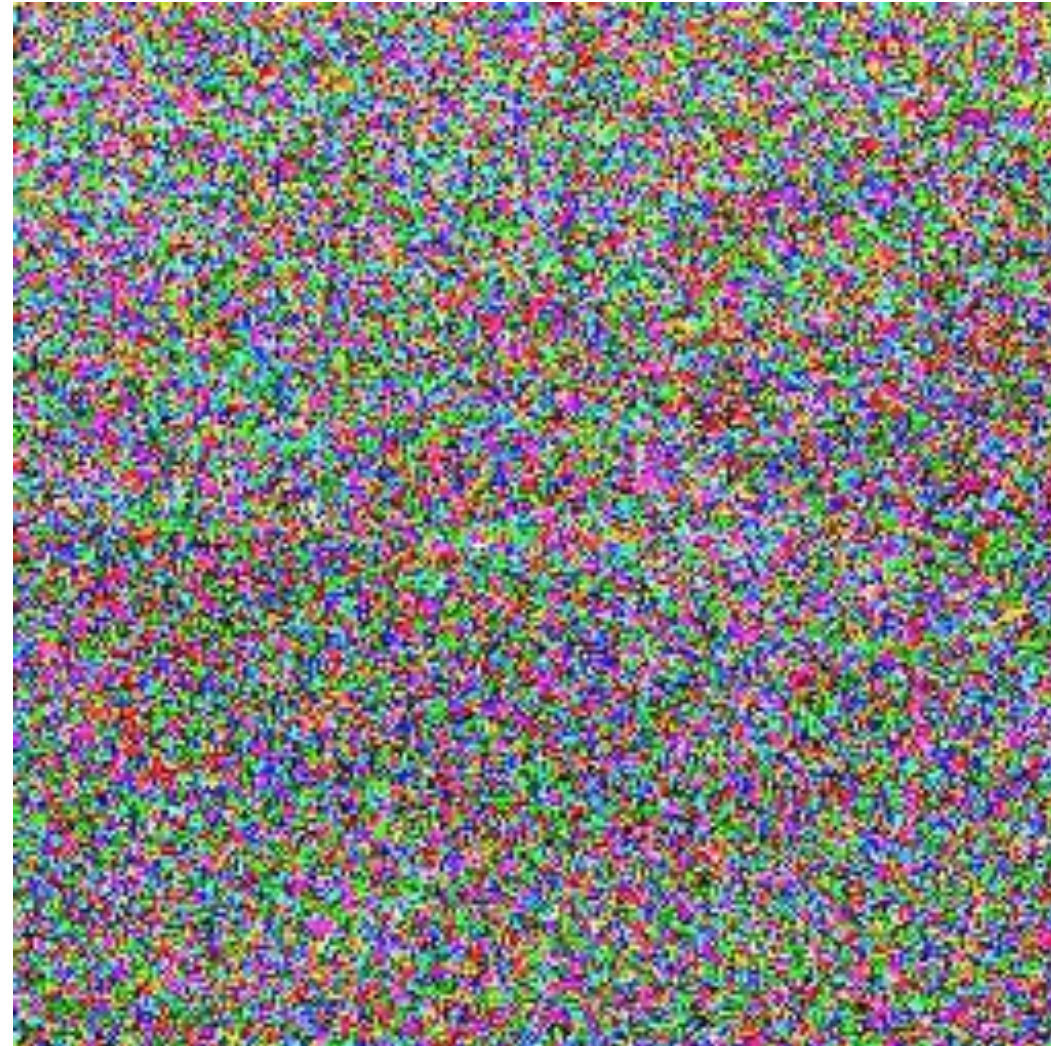


Figure 3. Multimodal behavior. At the given state, the end-effector (blue) can either go left or right to push the block. **Diffusion Policy** learns both modes and commits to only one mode within each rollout. In contrast, both **LSTM-GMM** **Mandlekar et al. (2021)** and **IBC** **Florence et al. (2021)** are biased toward one mode, while **BET** **Shafiullah et al. (2022)** fails to commit to a single mode due to its lack of temporal action consistency. Actions generated by rolling out 40 steps for the best-performing checkpoint.



2. High Dimensional Output Space



For a Image size of: $64*64$

Output Dimensions: 12288





2. High Dimensional Output Space

Allows joint inference of sequence of Actions.

X	807	33	274	458	-629	1062	-413	3	698	735	262	293	1603	339	492	571
Y	671	483	551	404	514	765	517	545	348	492	493	425	519	494	234	842

Toy example:

$(2 * 16) \rightarrow$ 32 dimensional output





2. High Dimensional Output Space

Allows joint inference of sequence of Actions.

X	807	33	274	458	-629	1062	-413	3	698	735	262	293	1603	339	492	571
Y	671	483	551	404	514	765	517	545	348	492	493	425	519	494	234	842

Toy example:

$(2 * 16) \rightarrow$ 32 dimensional output

6-DOF Manipulator:

$(6 * 16) \rightarrow$ 96 dimensional output





2. High Dimensional Output Space

Allows joint inference of sequence of Actions.

X	807	33	274	458	-629	1062	-413	3	698	735	262	293	1603	339	492	571
Y	671	483	551	404	514	765	517	545	348	492	493	425	519	494	234	842

Toy example:

$(2 * 16) \rightarrow$ 32 dimensional output

6-DOF Manipulator:

$(6 * 16) \rightarrow$ 96 dimensional output

Dual 6-DOF arm Manipulator:
with Gripper Control

$(14 * 16) \rightarrow$ 224 dimensional output



3. Stable Training

Because diffusion model directly predict the gradient of energy function



3. Stable Training

Because diffusion model directly predict the gradient of energy function

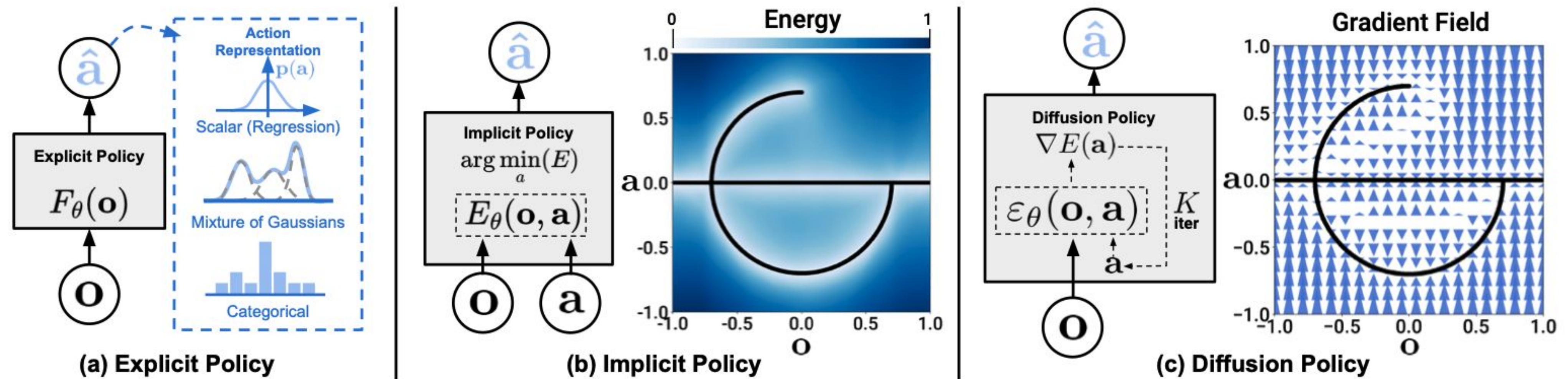


Figure 1. Policy Representations. a) Explicit policy with different types of action representations. b) Implicit policy learns an energy function conditioned on both action and observation and optimizes for actions that minimize the energy landscape c) Diffusion policy refines noise into actions via a learned gradient field. This formulation provides stable training, allows the learned policy to accurately model multimodal action distributions, and accommodates high-dimensional action sequences.



3. Stable Training

Because diffusion model directly predict the gradient of energy function

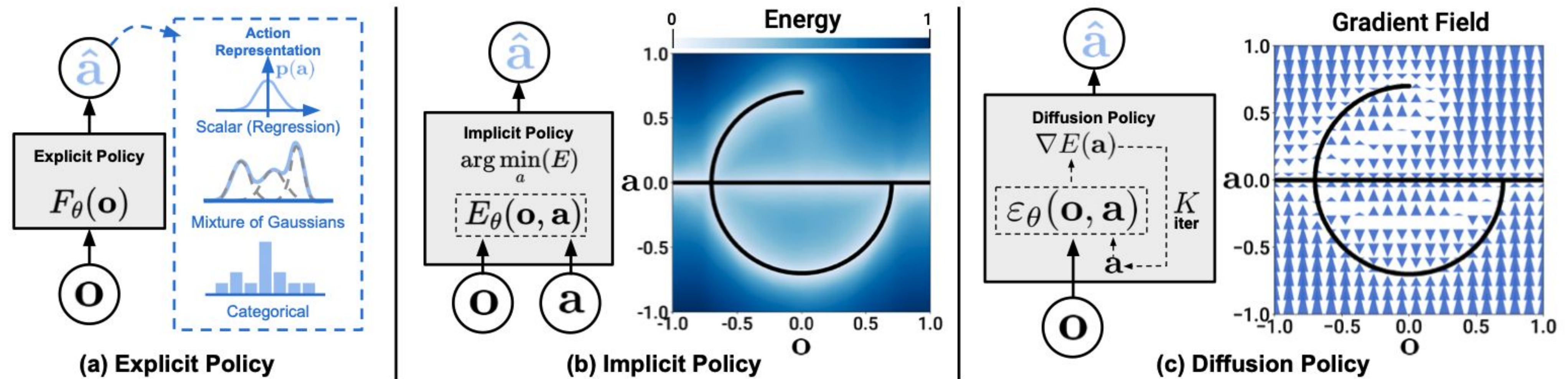


Figure 1. Policy Representations. a) Explicit policy with different types of action representations. b) Implicit policy learns an energy function conditioned on both action and observation and optimizes for actions that minimize the energy landscape c) Diffusion policy refines noise into actions via a learned gradient field. This formulation provides stable training, allows the learned policy to accurately model multimodal action distributions, and accommodates high-dimensional action sequences.



Small Maths Detour

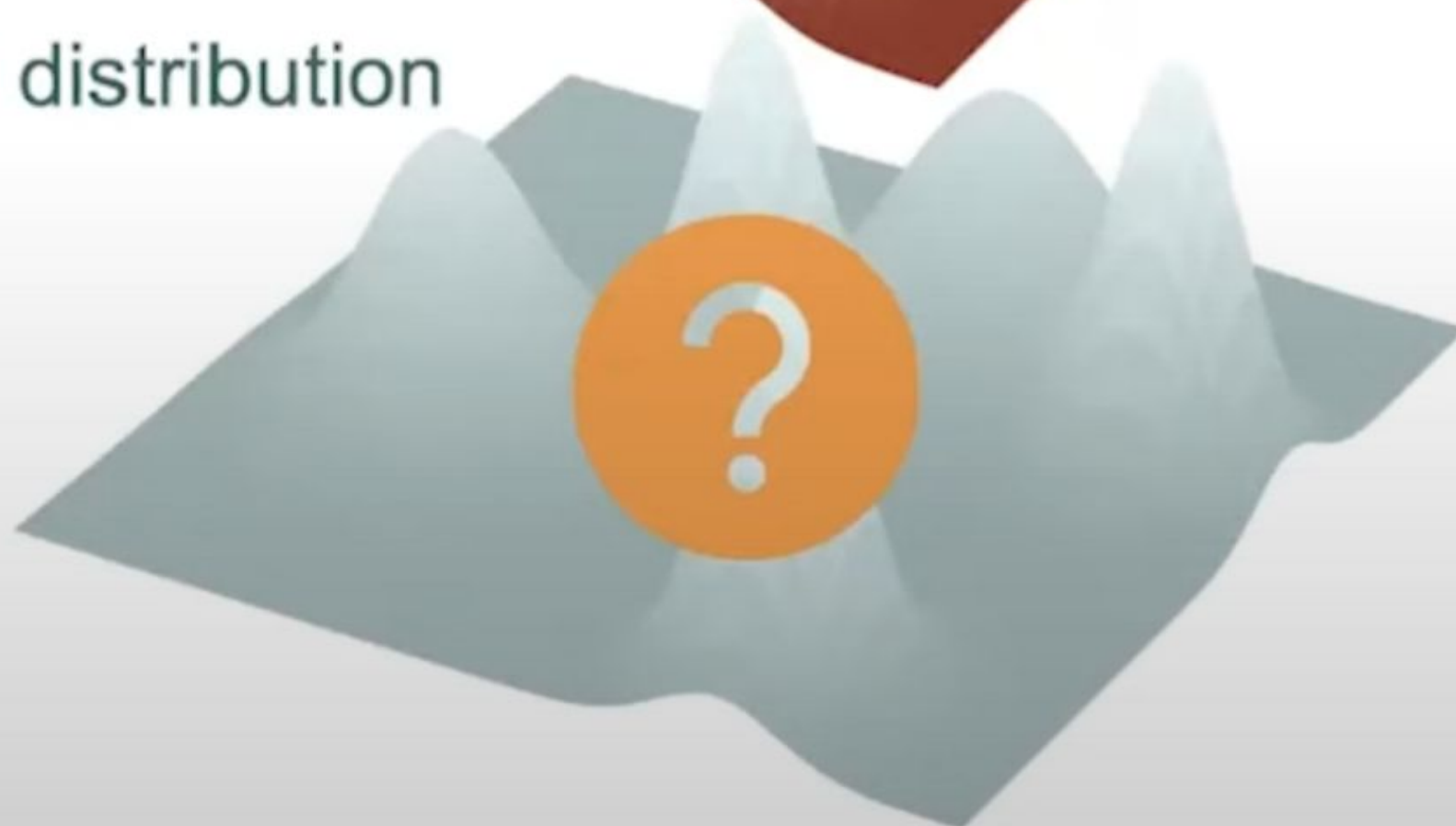
The key challenge for building complex generative models

Data distribution
(unknown)



Data distribution is extremely complex for high dimensional data.

Model distribution

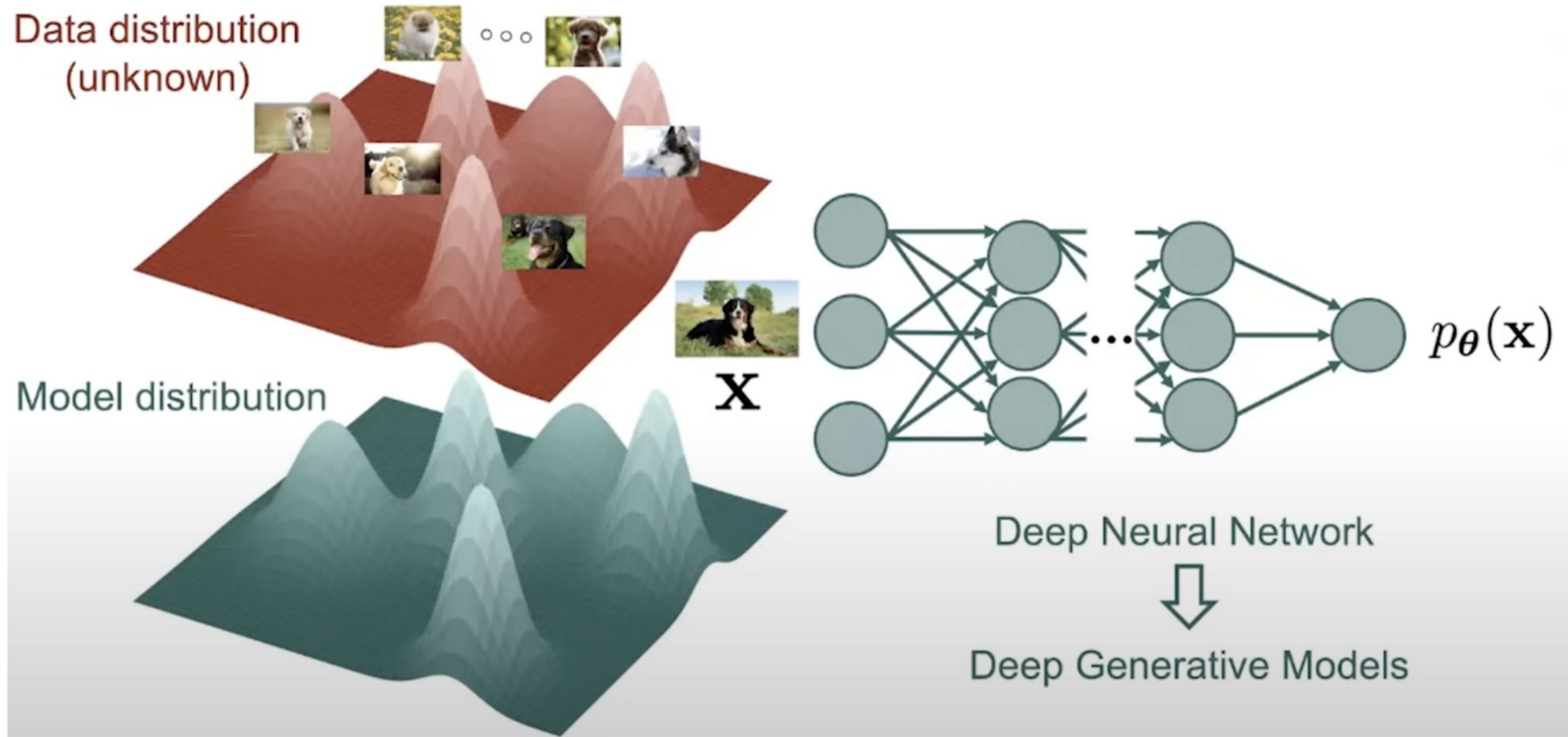


How to build a complex model to fit the data distribution?



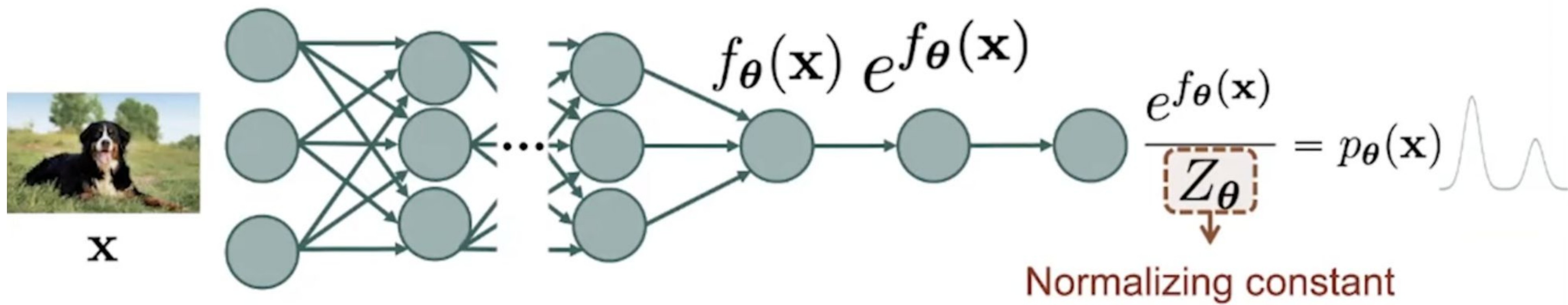
Small Maths Detour

The key challenge for building complex generative models



Small Maths Detour

The key challenge for building complex generative models



$$Z_{\theta} = \int e^{f_{\theta}(\mathbf{x})} d\mathbf{x}$$

X #P-complete even for discrete variables



Ludwig Boltzmann (1844-1906)



Willard Gibbs (1839-1903)



Gustav Zeuner (1828-1907)



Johannes van der Waals (1837-1923)

Thermodynamics & Statistical Mechanics



Small Maths Detour

Tackling the intractable normalizing constant

Approximating the normalizing constant

- Energy-based models [Ackley et al. 1985, LeCun et al. 2006]



Small Maths Detour

Tackling the intractable normalizing constant

Approximating the normalizing constant

- Energy-based models [Ackley et al. 1985, LeCun et al. 2006]



Inaccurate probability evaluation



Small Maths Detour

Tackling the intractable normalizing constant

Approximating the normalizing constant

- Energy-based models [Ackley et al. 1985, LeCun et al. 2006]



Inaccurate probability evaluation

Using restricted neural network models

- Autoregressive models [Bengio & Bengio 2000, van den Oord et al. 2016]
- Normalizing flow models [Dinh et al. 2014, Rezende & Mohamed 2015]
- Variational auto-encoders [Kingma & Welling 2014, Rezende et al. 2014]





Small Maths Detour

Tackling the intractable normalizing constant

Approximating the normalizing constant

- Energy-based models [Ackley et al. 1985, LeCun et al. 2006]



Inaccurate probability evaluation

Using restricted neural network models

- Autoregressive models [Bengio & Bengio 2000, van den Oord et al. 2016]
- Normalizing flow models [Dinh et al. 2014, Rezende & Mohamed 2015]
- Variational auto-encoders [Kingma & Welling 2014, Rezende et al. 2014]



Restricted model family



Small Maths Detour

Tackling the intractable normalizing constant

Approximating the normalizing constant

- Energy-based models [Ackley et al. 1985, LeCun et al. 2006]



Inaccurate probability evaluation

Using restricted neural network models

- Autoregressive models [Bengio & Bengio 2000, van den Oord et al. 2016]
- Normalizing flow models [Dinh et al. 2014, Rezende & Mohamed 2015]
- Variational auto-encoders [Kingma & Welling 2014, Rezende et al. 2014]



Restricted model family

Modeling the Generation Process Only

- Generative Adversarial Networks (GANs) [Goodfellow et al. 2014]





Small Maths Detour

Tackling the intractable normalizing constant

Approximating the normalizing constant

- Energy-based models [Ackley et al. 1985, LeCun et al. 2006]



Inaccurate probability evaluation

Using restricted neural network models

- Autoregressive models [Bengio & Bengio 2000, van den Oord et al. 2016]
- Normalizing flow models [Dinh et al. 2014, Rezende & Mohamed 2015]
- Variational auto-encoders [Kingma & Welling 2014, Rezende et al. 2014]



Restricted model family

Modeling the Generation Process Only

- Generative Adversarial Networks (GANs) [Goodfellow et al. 2014]

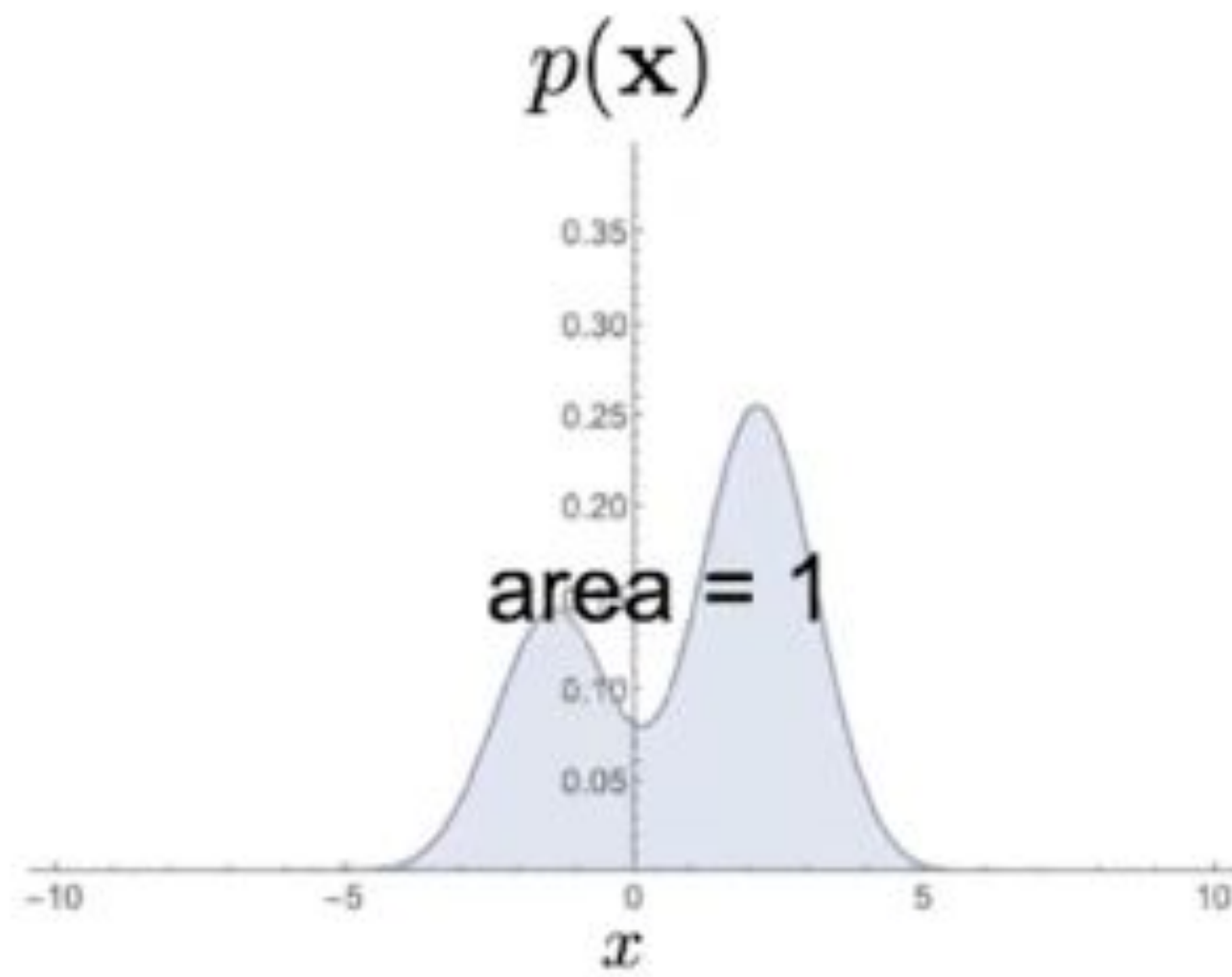


Cannot evaluate probabilities

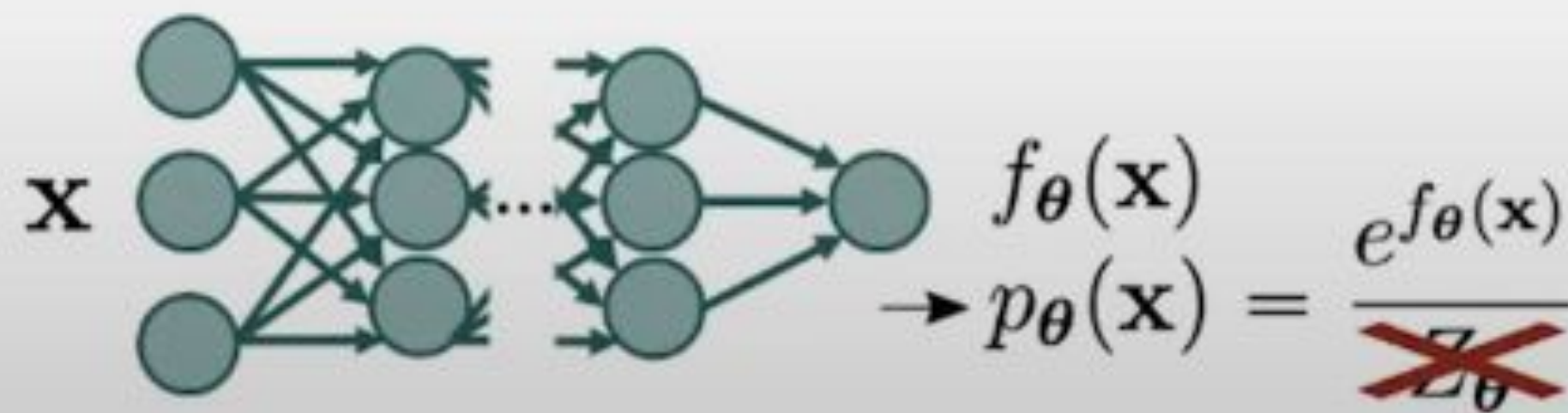




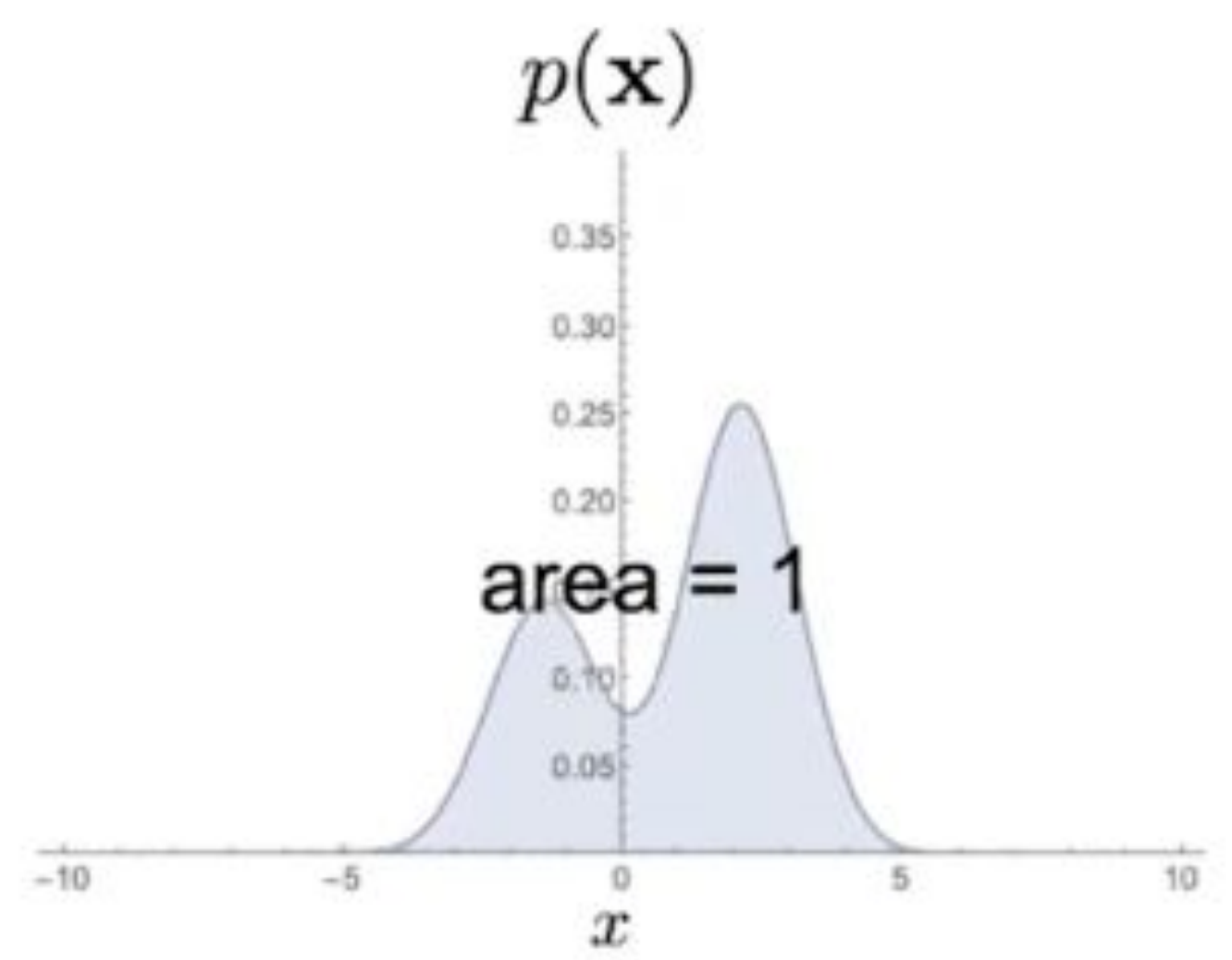
Small Maths Detour



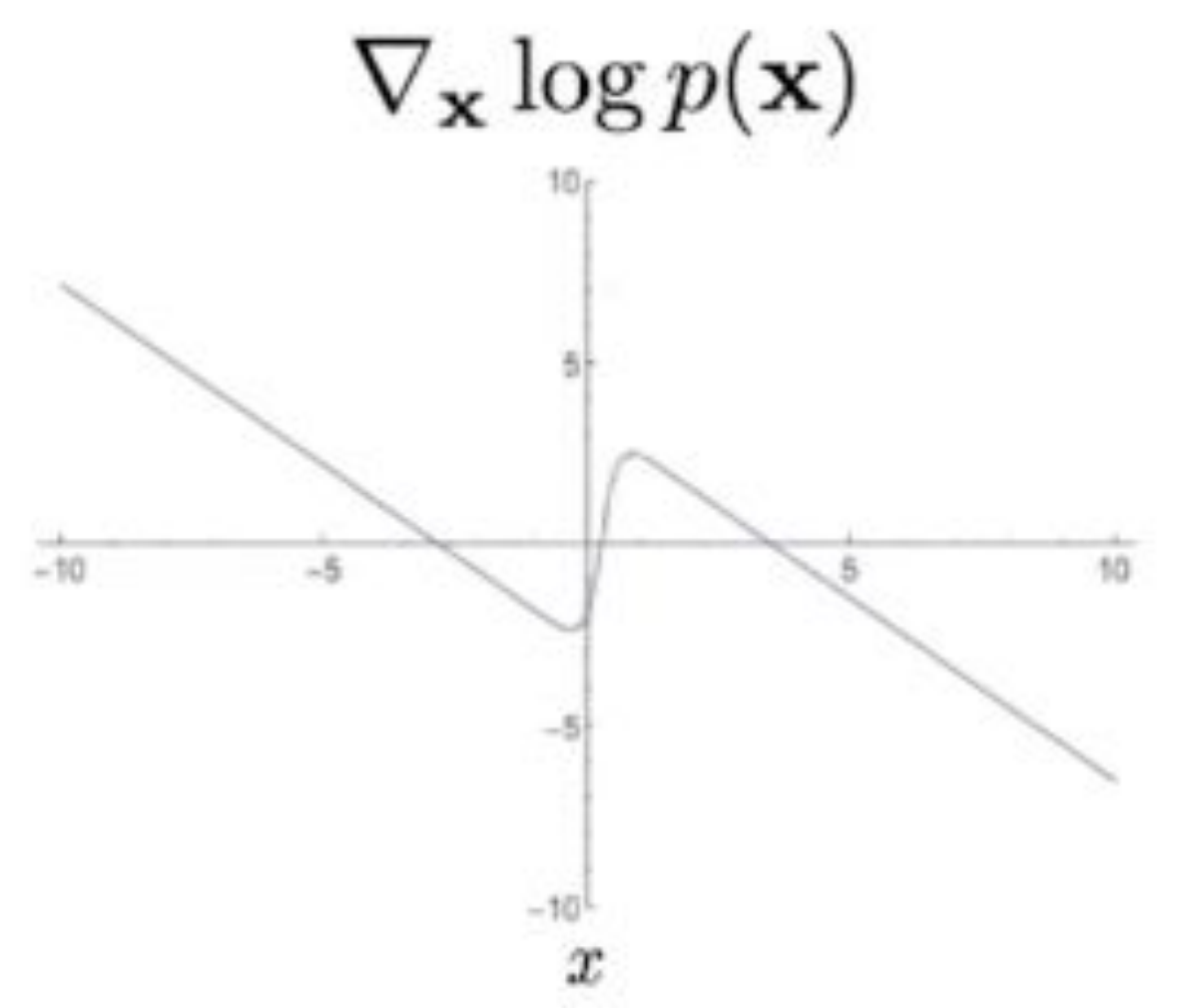
Probability density function



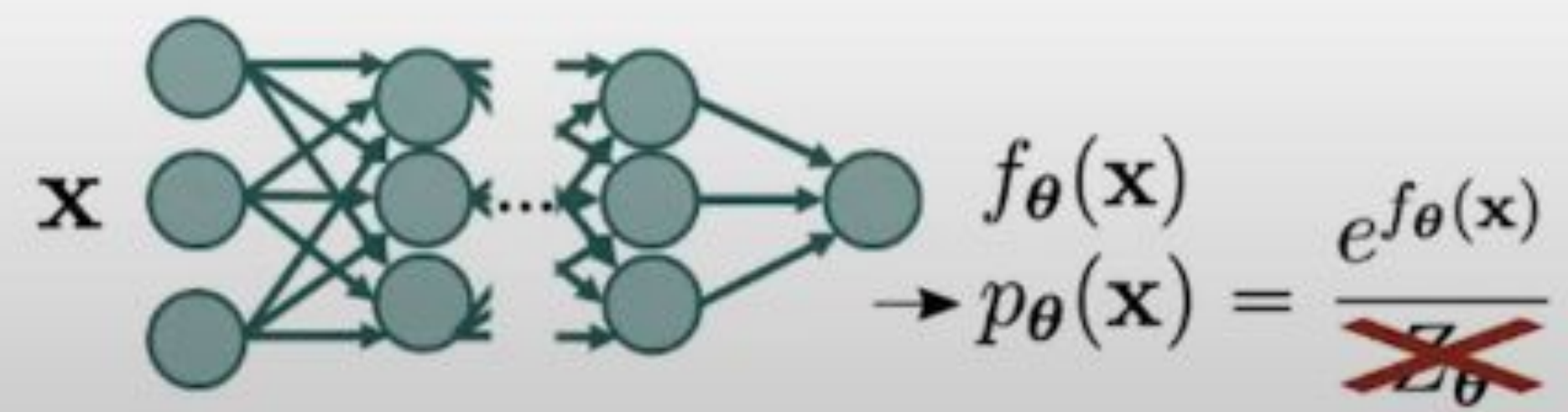
Small Maths Detour



Probability density function



Score function



$$\begin{aligned}
 \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}) &= \nabla_{\mathbf{x}} f_{\theta}(\mathbf{x}) - \nabla_{\mathbf{x}} \log Z_{\theta} \\
 &= \nabla_{\mathbf{x}} f_{\theta}(\mathbf{x}) \\
 &= \mathbf{s}_{\theta}(\mathbf{x})
 \end{aligned}$$

$\nabla_{\mathbf{x}} \log Z_{\theta} \downarrow 0$



3. Stable Training

Because diffusion model directly predict the gradient of energy function

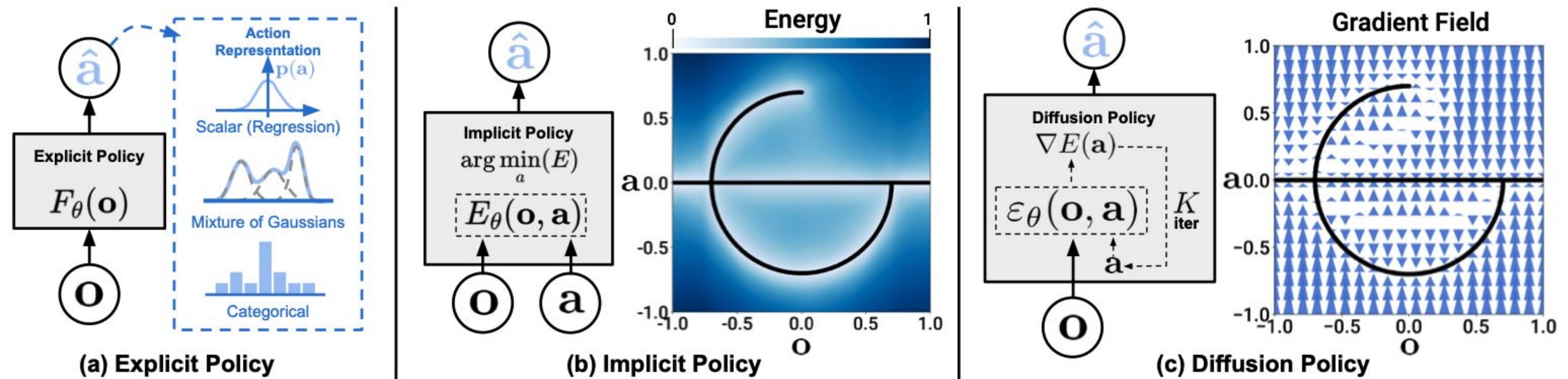


Figure 1. Policy Representations. a) Explicit policy with different types of action representations. b) Implicit policy learns an energy function conditioned on both action and observation and optimizes for actions that minimize the energy landscape c) Diffusion policy refines noise into actions via a learned gradient field. This formulation provides stable training, allows the learned policy to accurately model multimodal action distributions, and accommodates high-dimensional action sequences.





3. Stable Training

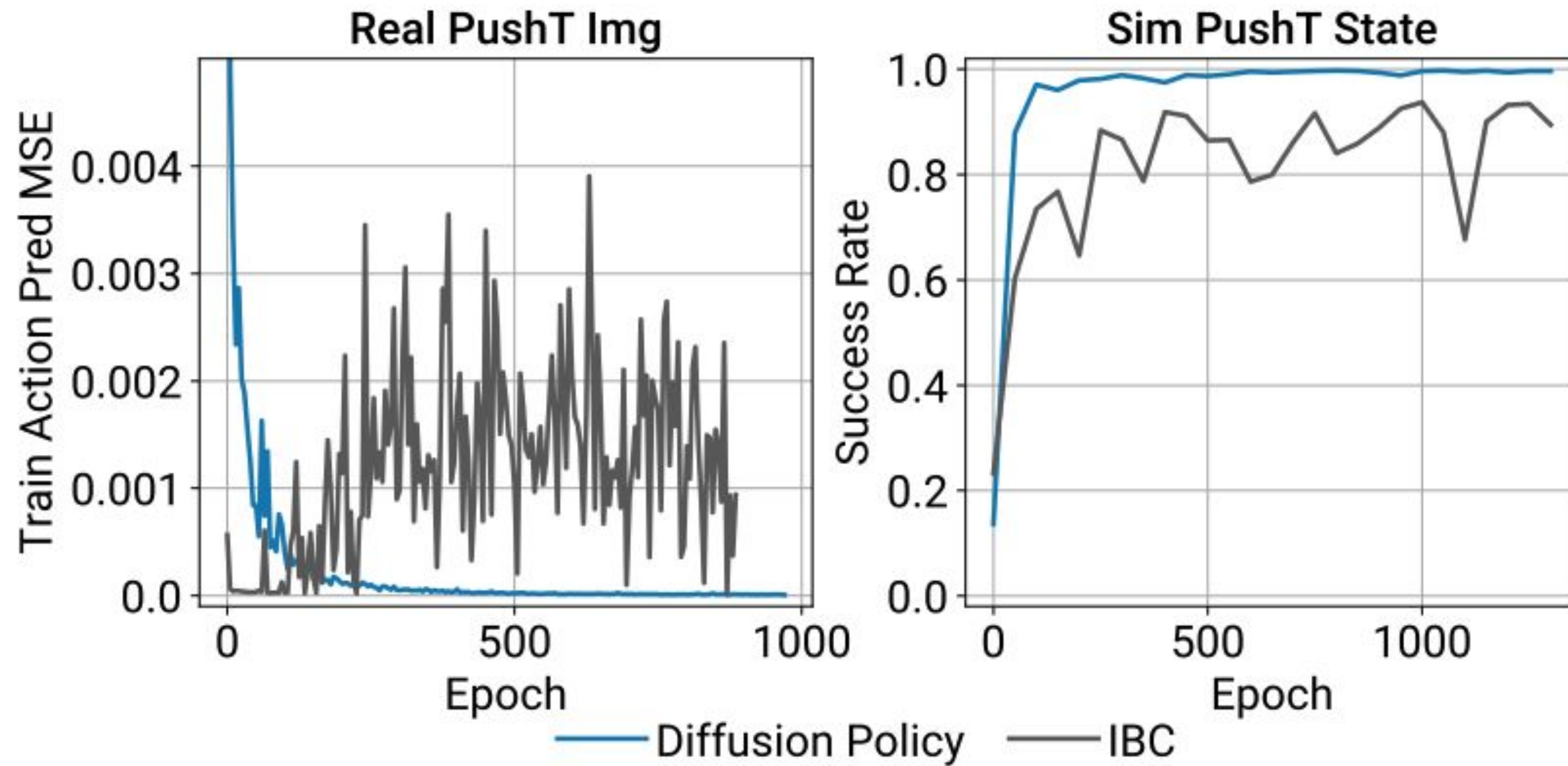


Figure 6. Training Stability. Left: IBC fails to infer training actions with increasing accuracy despite smoothly decreasing training loss for energy function. Right: IBC’s evaluation success rate oscillates, making checkpoint selection difficult (evaluated using policy rollouts in simulation).



Overview

- Motivation behind diffusion models ✓
- Denoising Diffusion Probabilistic Models ✓
 - Diffusion Process and Reverse Diffusion Process ✓
 - Training and Sampling a Diffusion Model ✓
 - Results from a Diffusion Model that we trained ✓
- Conditioned v/s Unconditioned Diffusion ✓
- Policy Learning ✓
- Toy Problem ✓
- Diffusion Policy on Real Robots
 - Advantages ✓
 - Network Architect
 - Evaluation & Results





Closed Loop Action-Sequence Prediction





Closed Loop Action-Sequence Prediction

t

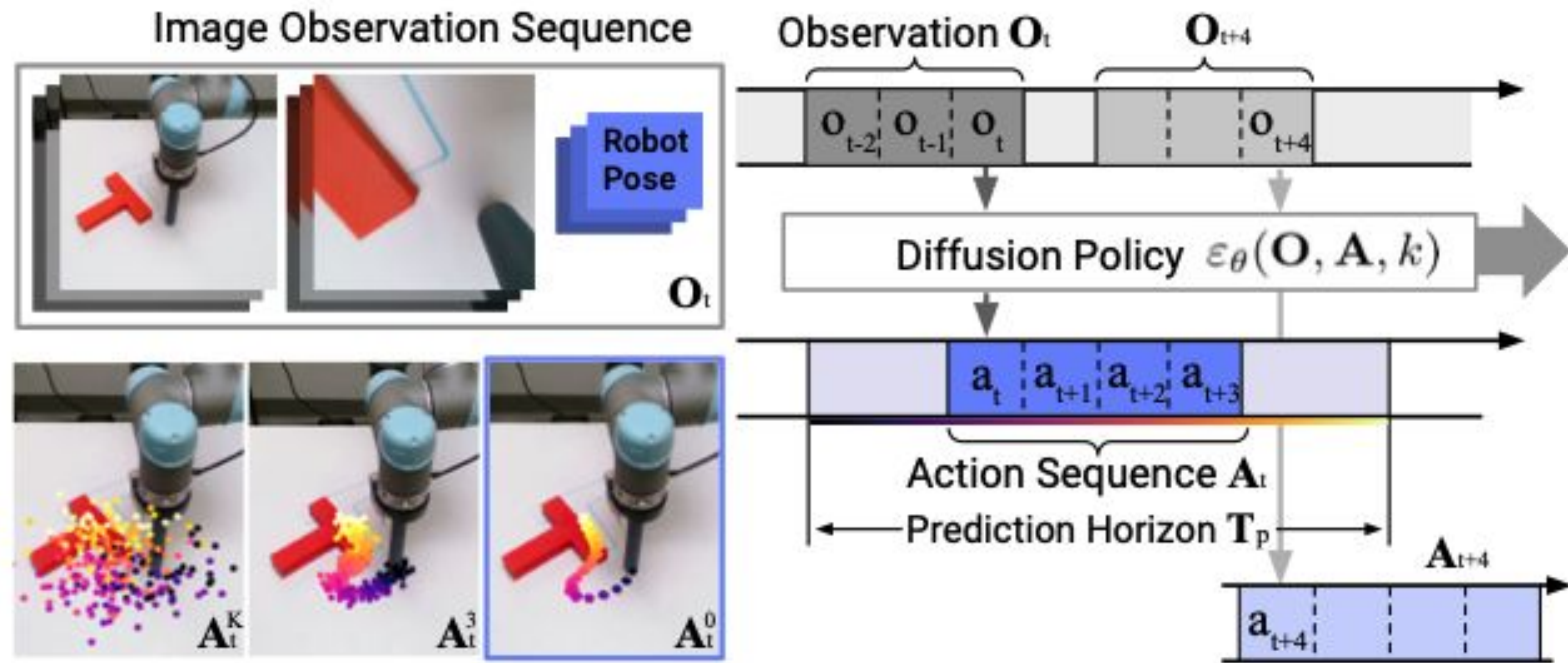
Observation:																			
Action Pred:																			
Action Used:																			
time:	-4	-3	-2	-1	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14





Closed Loop Action-Sequence Prediction

	t																		
Observation:																			
Action Pred:																			
Action Used:																			
time:	-4	-3	-2	-1	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14



Output: Action Sequence

a) Diffusion Policy General Formulation



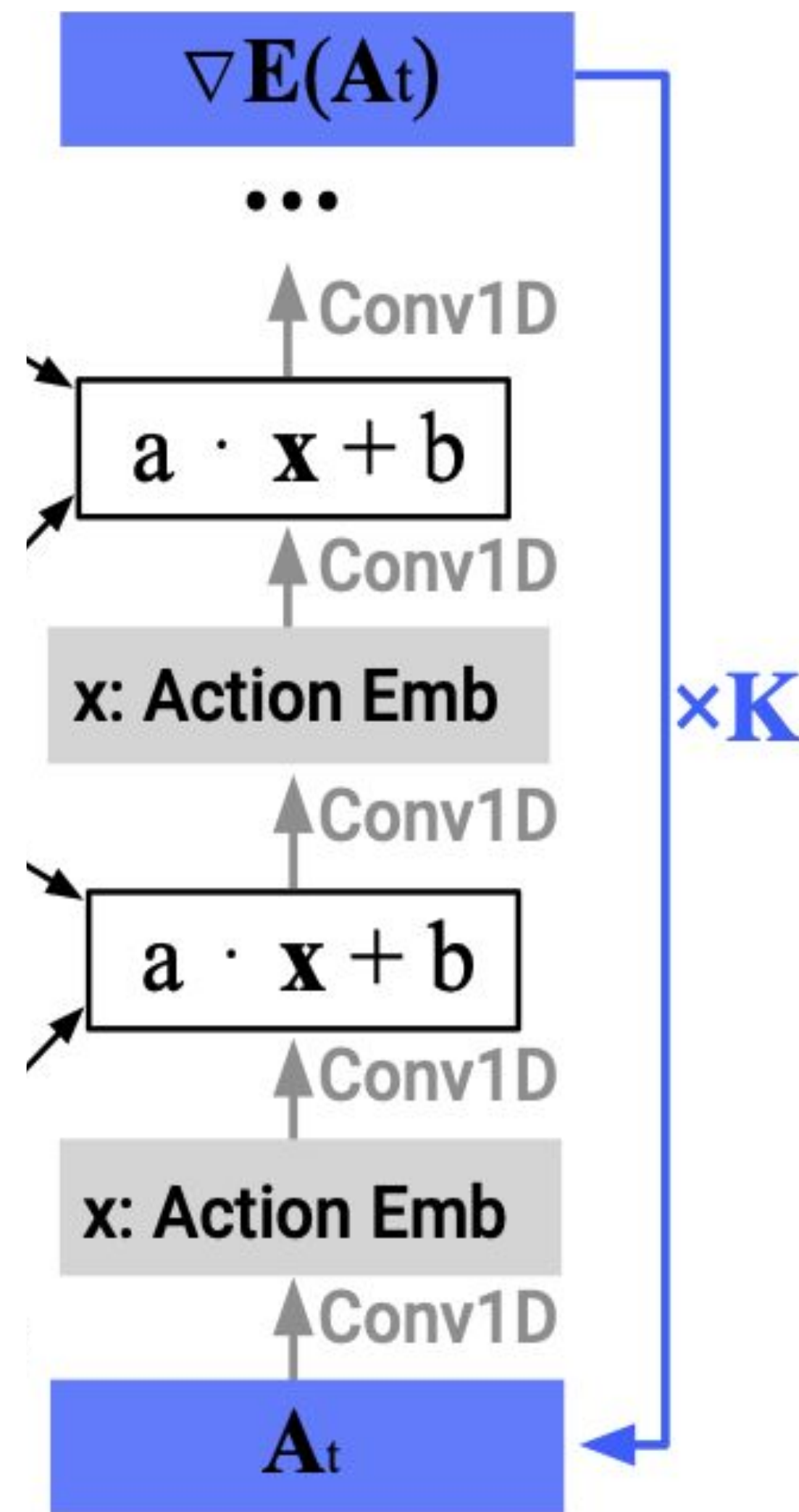


Network Architecture for ϵ_θ





Network Architecture for \mathcal{E}_θ

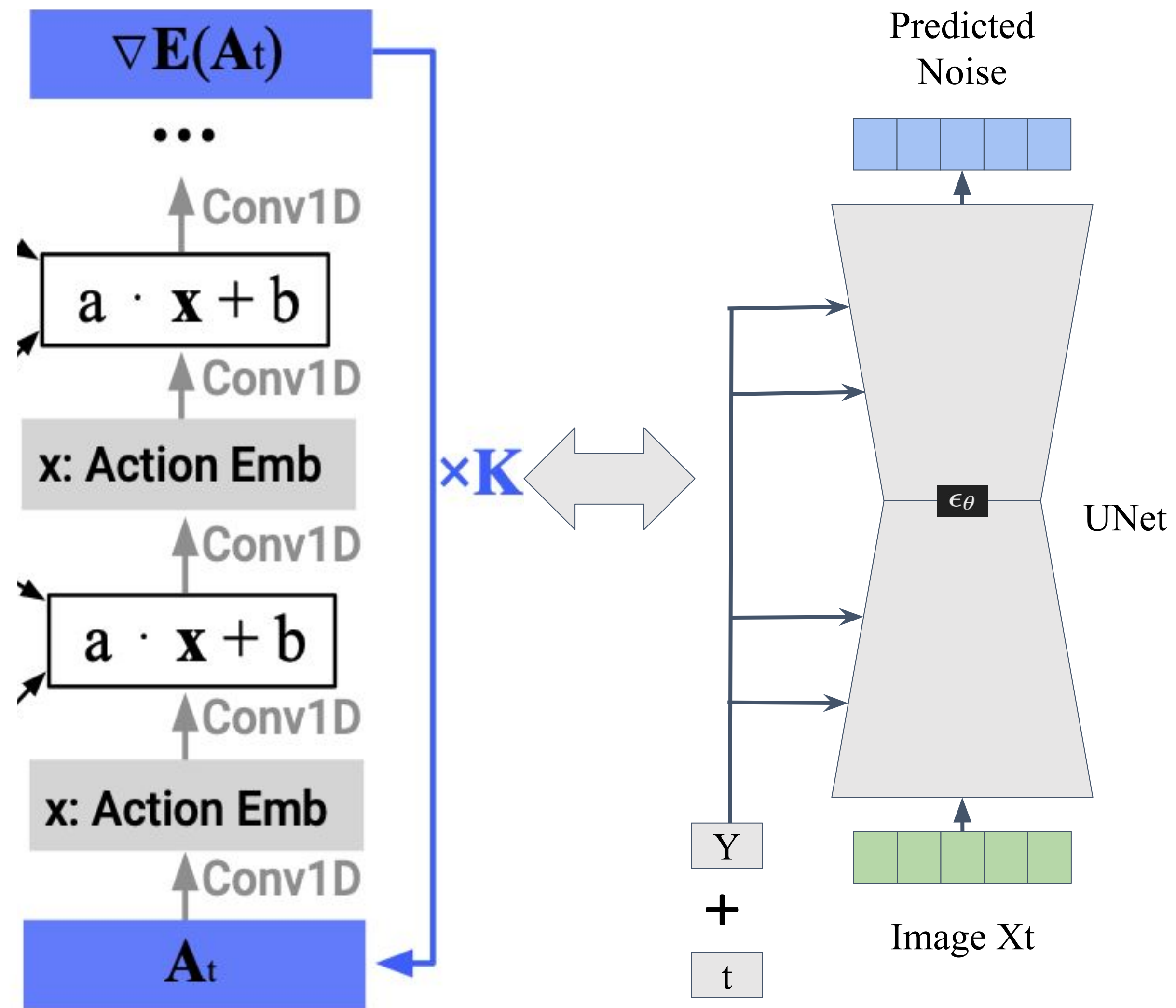


C. Chi *et al.*, "Diffusion Policy: Visuomotor Policy Learning via Action Diffusion," *arXiv preprint arXiv:2303.04137*, 2024. [Online]. Available: <https://arxiv.org/abs/2303.04137>

Image Source: <https://www.assemblyai.com/blog/how-imagen-actually-works/>

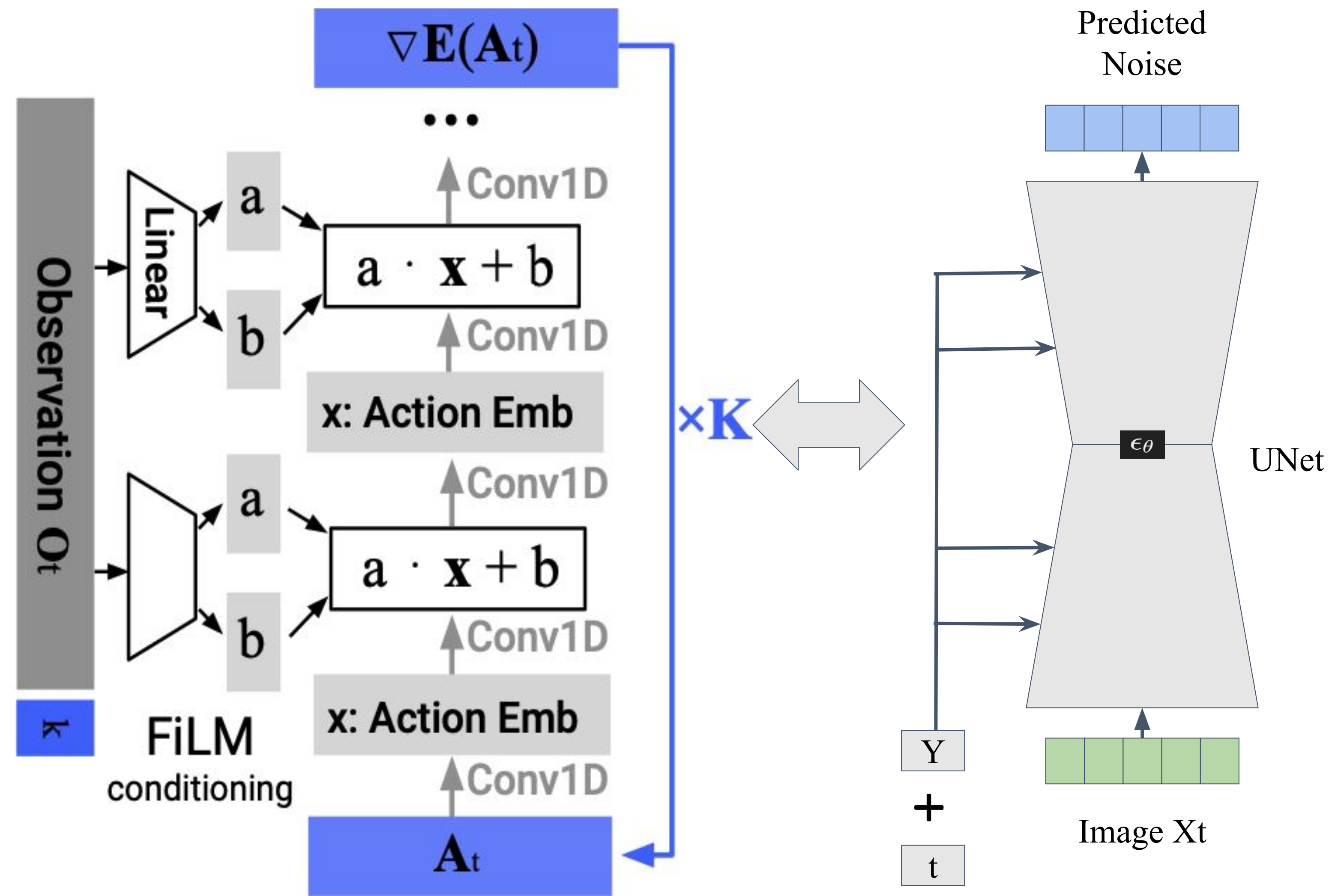


Network Architecture for ϵ_θ





Network Architecture for ϵ_θ

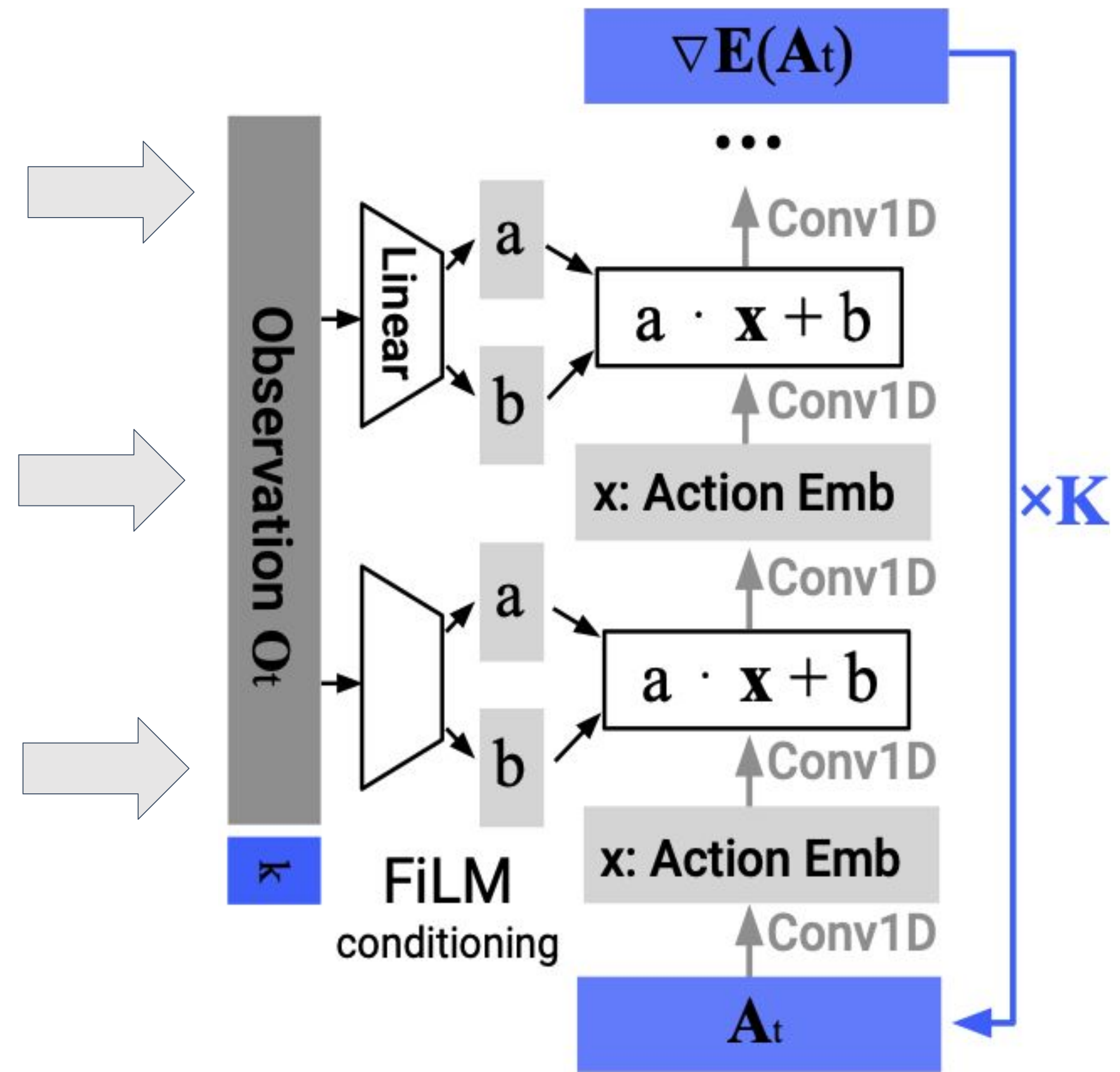


C. Chi *et al.*, "Diffusion Policy: Visuomotor Policy Learning via Action Diffusion," *arXiv preprint arXiv:2303.04137*, 2024. [Online]. Available: <https://arxiv.org/abs/2303.04137>

Image Source: <https://www.assemblyai.com/blog/how-imagen-actually-works/>

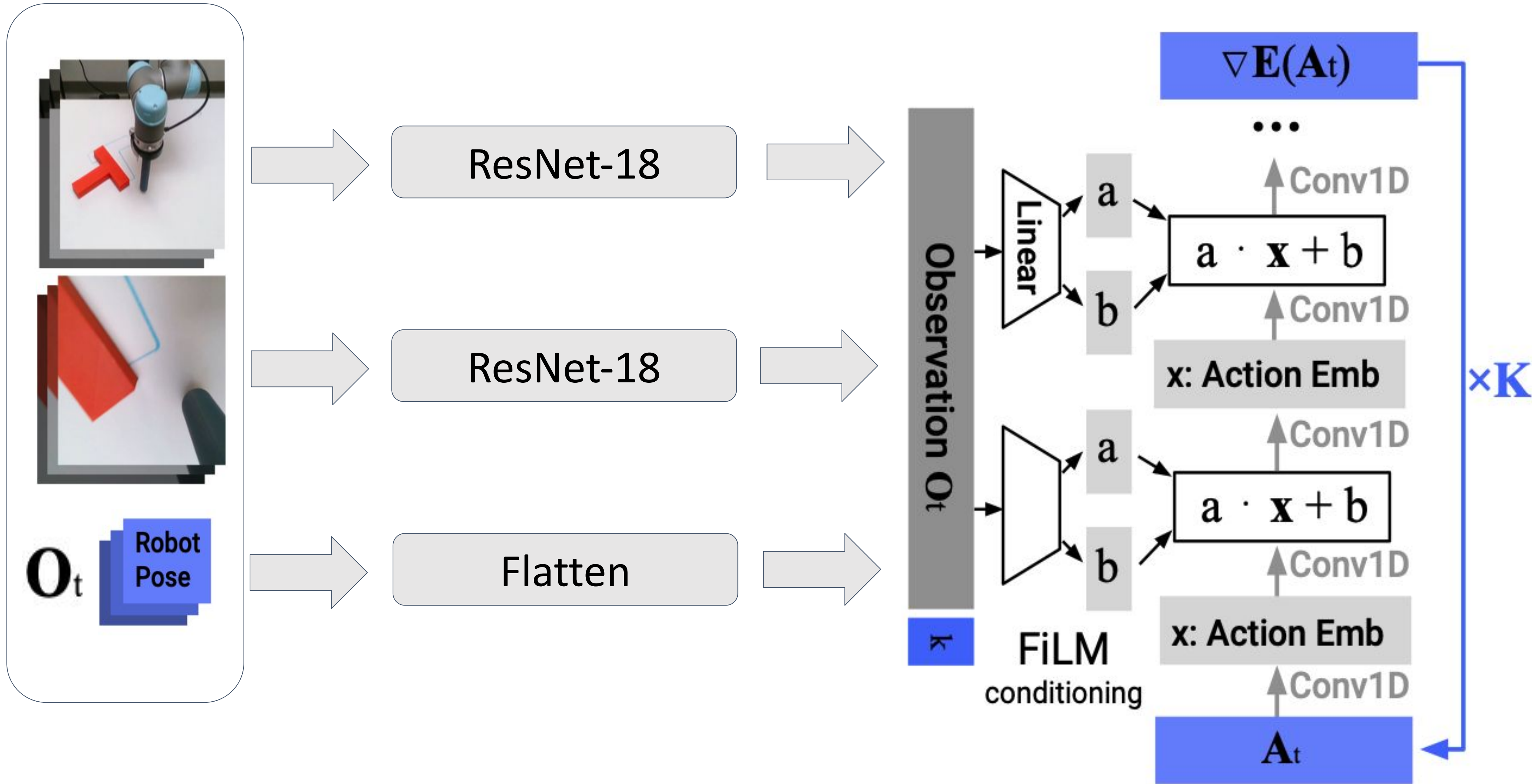


Network Architecture for \mathcal{E}_θ





Network Architecture for ϵ_θ



Overview

- Motivation behind diffusion models ✓
- Denoising Diffusion Probabilistic Models ✓
 - Diffusion Process and Reverse Diffusion Process ✓
 - Training and Sampling a Diffusion Model ✓
 - Results from a Diffusion Model that we trained ✓
- Conditioned v/s Unconditioned Diffusion ✓
- Policy Learning ✓
- Toy Problem ✓
- Diffusion Policy on Real Robots
 - Advantages ✓
 - Network Architect ✓
 - Evaluation & Results



Evaluation

Real World Push-T task:

- trained with 136 Proficient-Human demonstrations

	Human	IBC		LSTM-GMM		Diffusion Policy			
	Demo	pos	vel	pos	vel	T-E2E	ImgNet	R3M	E2E
IoU	0.84	0.14	0.19	0.24	0.25	0.53	0.24	0.66	0.80
Succ%	1.00	0.00	0.00	0.20	0.10	0.65	0.15	0.80	0.95
Dur.	20.3	56.3	41.6	47.3	51.7	57.5	55.8	31.7	22.9

- Results on 20 rollouts
- Rollout is Successful if $\text{IoU}_{\text{Last Step}} > \min(\text{IoU from human demonstrations})$
- Max. steps 600 in a rollout





Evaluation

Real World Push-T task:

- trained with 136 Proficient-Human demonstrations





Visual Results

Let's go to the website to see the policies performance in action:

<https://diffusion-policy.cs.columbia.edu>



Overview

- Motivation behind diffusion models ✓
- Denoising Diffusion Probabilistic Models ✓
 - Diffusion Process and Reverse Diffusion Process ✓
 - Training and Sampling a Diffusion Model ✓
 - Results from a Diffusion Model that we trained ✓
- Conditioned v/s Unconditioned Diffusion ✓
- Policy Learning ✓
- Toy Problem ✓
- Diffusion Policy on Real Robots ✓
 - Advantages ✓
 - Network Architect ✓
 - Evaluation & Results ✓



Next Lecture:

Student Lecture 7

Dual-arm Manipulation





Happy Thanksgiving!
No class on 11/27



DR

DeepRob

[Group 8] Lecture 6

Diffusion Models and Policy Learning

by *Nirshal Chandra Sekar and Mohit Yadav*

University of Minnesota



DR

DeepRob

[Group 8] Lecture 6

Diffusion Models and Policy Learning

by *Nirshal Chandra Sekar and Mohit Yadav*

University of Minnesota

