# Imitation Learning Tutorial

Chenyu Yang
Davide Liconti
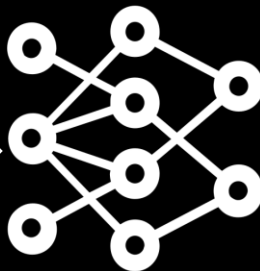
How does what we learned so far connect to each other?

Design concepts, CAD
Unit 2

Fabrication
Unit 3

Kinematics & Control
Unit 4

Robot arm & Dynamixel controller

STL meshes,
joint structure

Simulation
Unit 4

Physical
robot hand

Teleoperation
Unit 6

MJCF simulation model

RL / IL
Unit 7

Applications

ETH zürich    SoftRobotics Laboratory

# Overview

**Observations**

**Actions**

$\mathcal{Z}_{obs} \rightarrow \mathcal{Z}_{act}$

What observations can we use?
What's the best $\mathcal{Z}_{obs}$?
What's the best $\mathcal{Z}_{act}$?
What's the best way to learn?

For robotic hands !
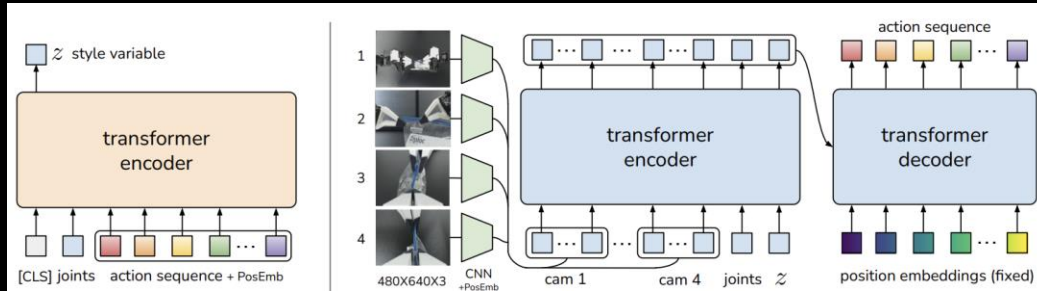
ETHzürich    SoftRobotics Laboratory

# Contents

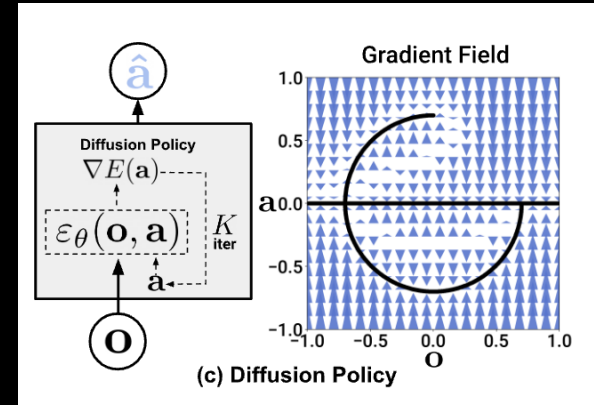2 architectures for downstream-task imitation learning for robotics
ACT ([paper](#))
Diffusion Policy([paper](#))

What we need to understand today:
1. Transformer Architecture
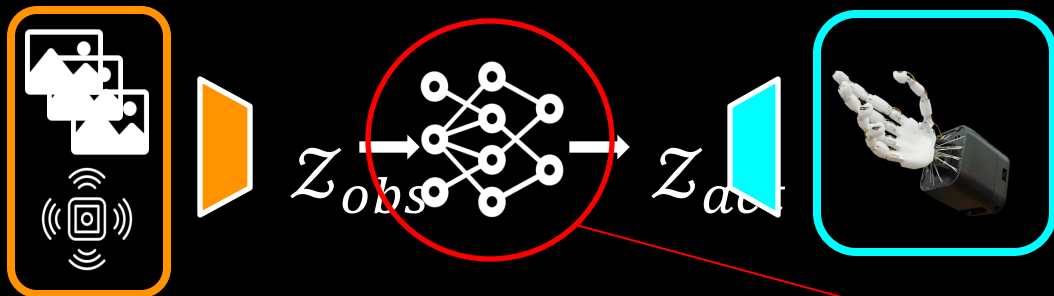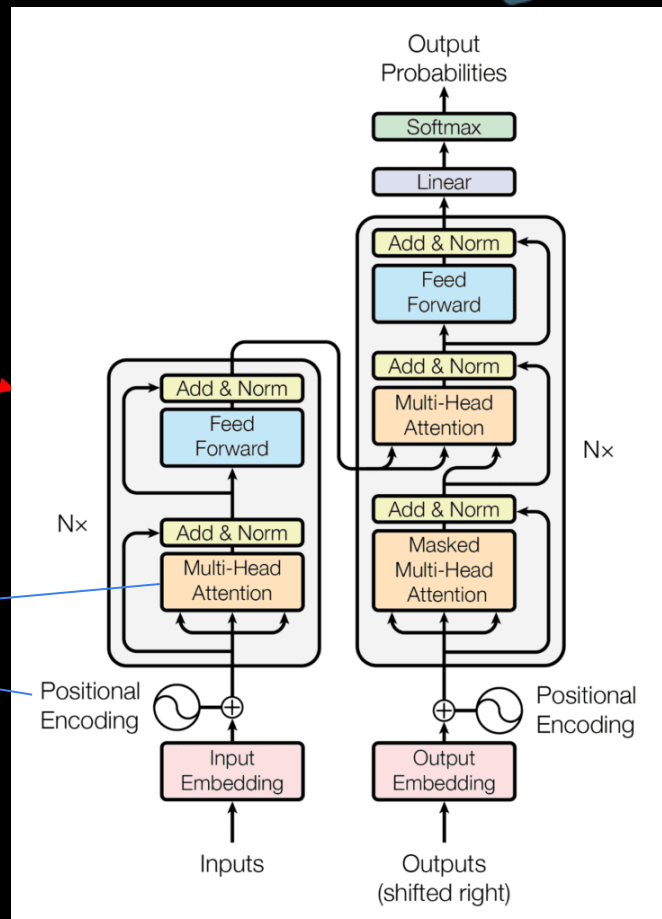2. Variational Autoencoder
3. Diffusion Mechanism



ACT



Diffusion Policy

# Transformers



$z_{obs}$ → $z_{act}$

The core model is usually a transformer-based architecture
(ACT, OCTO, OpenVLA, …)

1. Positional Encoding
2. Self-Attention



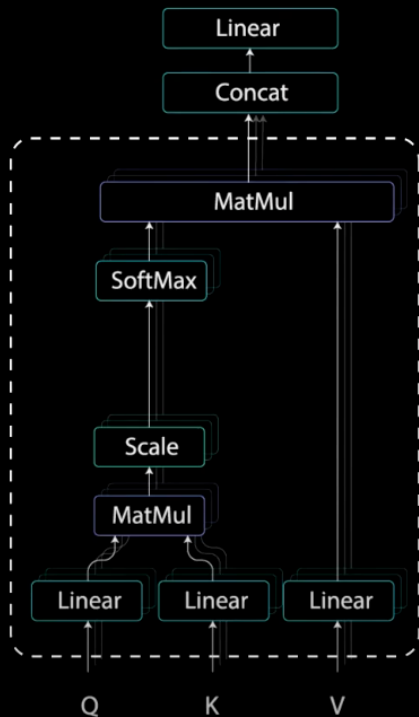**ETH**zürich   **S**oftRobotics Laboratory

# Positional Embeddings

- Unlike RNNs, transformers don't inherently understand sequence order Positional embeddings add information about token positions in sequences.

- Positional information is added to token embeddings, often using sinusoidal functions to allow model generalization to various sequence lengths.

$$PE_{pos,2i} = \sin\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right) \quad PE_{pos,2i+1} = \cos\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$$
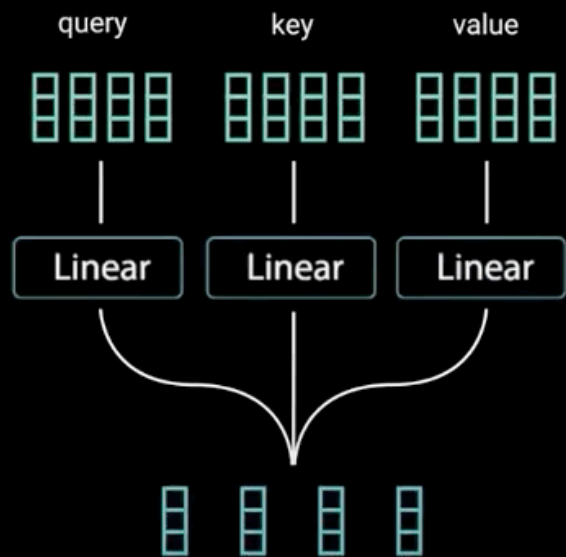
# Attention



From input embeddinsg use 3 distinct linear layers to create query, key and value vectors

"For example, when you type a query to search for some video on Youtube, the search engine will map your **query** against a set of **keys** (video title, description etc.) associated with candidate videos in the database, then present you the best matched videos (**values**)."

# Attention

# Attention

Softmax( ⊞ ) =

|  | Hi | how | are | you |
|-----|-----|-----|-----|-----|
| Hi | 0.7 | 0.1 | 0.1 | 0.1 |
| how | 0.1 | 0.6 | 0.2 | 0.1 |
| are | 0.1 | 0.3 | 0.6 | 0.1 |
| you | 0.1 | 0.3 | 0.3 | 0.3 |

$$softmax(x)_i = \frac{exp(x_i)}{\sum_j exp(x_j))}$$

attention weights        value        output

 X  = 

ETH zürich        SoftRobotics Laboratory

# Decoder

The decoder is autoregressive
-> How to prevent conditioning on future tokens?



HI,

how

are

you?

Transformers
Decoder

<start>

|       | <start> | I | am | fine |
|-------|---------|-----|-----|------|
| <start> | 0.7 | 0.1 | 0.1 | 0.1 |
| I | 0.1 | 0.6 | 0.2 | 0.1 |
| am | 0.1 | 0.3 | 0.6 | 0.1 |
| fine | 0.1 | 0.3 | 0.3 | 0.3 |

Scaled Scores

| 0.7 | 0.1 | 0.1 | 0.1 |
| 0.1 | 0.6 | 0.2 | 0.1 |
| 0.1 | 0.3 | 0.6 | 0.1 |
| 0.1 | 0.3 | 0.3 | 0.3 |

**+**

Look-Ahead Mask

| 0 | -inf | -inf | -inf |
| 0 | 0 | -inf | -inf |
| 0 | 0 | 0 | -inf |
| 0 | 0 | 0 | 0 |

**=**

Masked Scores

| 0.7 | -inf | -inf | -inf |
| 0.1 | 0.6 | -inf | -inf |
| 0.1 | 0.3 | 0.6 | -inf |
| 0.1 | 0.3 | 0.3 | 0.3 |

ETH zürich

SoftRobotics Laboratory

# Transformer explainer

- [https://poloclub.github.io/transformer-explainer/](https://poloclub.github.io/transformer-explainer/)

# Variational Autoencoder

- Autoencoder

Input $x$



Encoder

Bottleneck
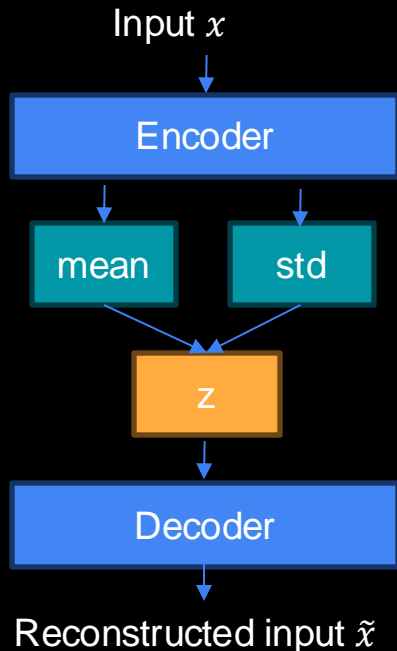
Decoder

Reconstructed input $\tilde{x}$

```python
def forward(self, x):
    z = self.encode(x)
    return self.decode(z)

def loss_function(recon_x, x):
    l2 = F.mse_loss(recon_x, x, reduction='sum')
    return l2
```

ETH zürich   SoftRobotics Laboratory

# Variational Autoencoder

- Variational Autoencoder

Input $x$



Encoder

mean | std

z

Decoder

Reconstructed input $\tilde{x}$

```python
def reparameterize(self, mu, logvar):
    std = torch.exp(0.5 * logvar)
    eps = torch.randn_like(std)
    return mu + eps * std

def forward(self, x):
    mu, logvar = self.encode(x)
    z = self.reparameterize(mu, logvar)
    return self.decode(z), mu, logvar

def loss_function(recon_x, x, mu, logvar):
    kld = -0.5 * torch.sum(1 + logvar - mu.pow(2) - logvar.exp())
    BCE = F.mse_loss(recon_x, x, reduction='sum')
    return BCE
```

Link: **INTERPRETING LATENT SPACE AND BIAS**

# Conditional Variational Autoencoder

- Variational Autoencoder



```
def reparameterize(self, mu, logvar):
    std = torch.exp(0.5 * logvar)
    eps = torch.randn_like(std)
    return mu + eps * std

def forward(self, x, obs):
    mu, logvar = self.encode(x)
    z = self.reparameterize(mu, logvar)
    return self.decode(z, self.obs_encode(obs))

def loss_function(recon_x, x, mu, logvar):
    kld = -0.5 * torch.sum(1 + logvar - mu.pow(2) - logvar.exp())
    BCE = F.mse_loss(recon_x, x, reduction='sum')
    return BCE
```
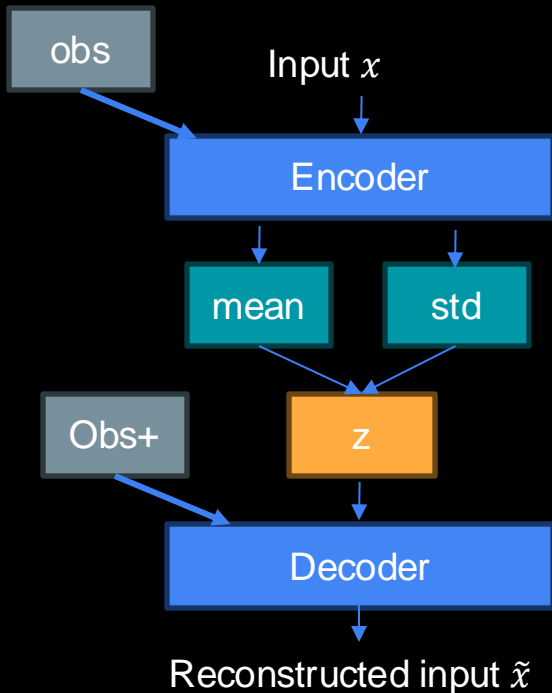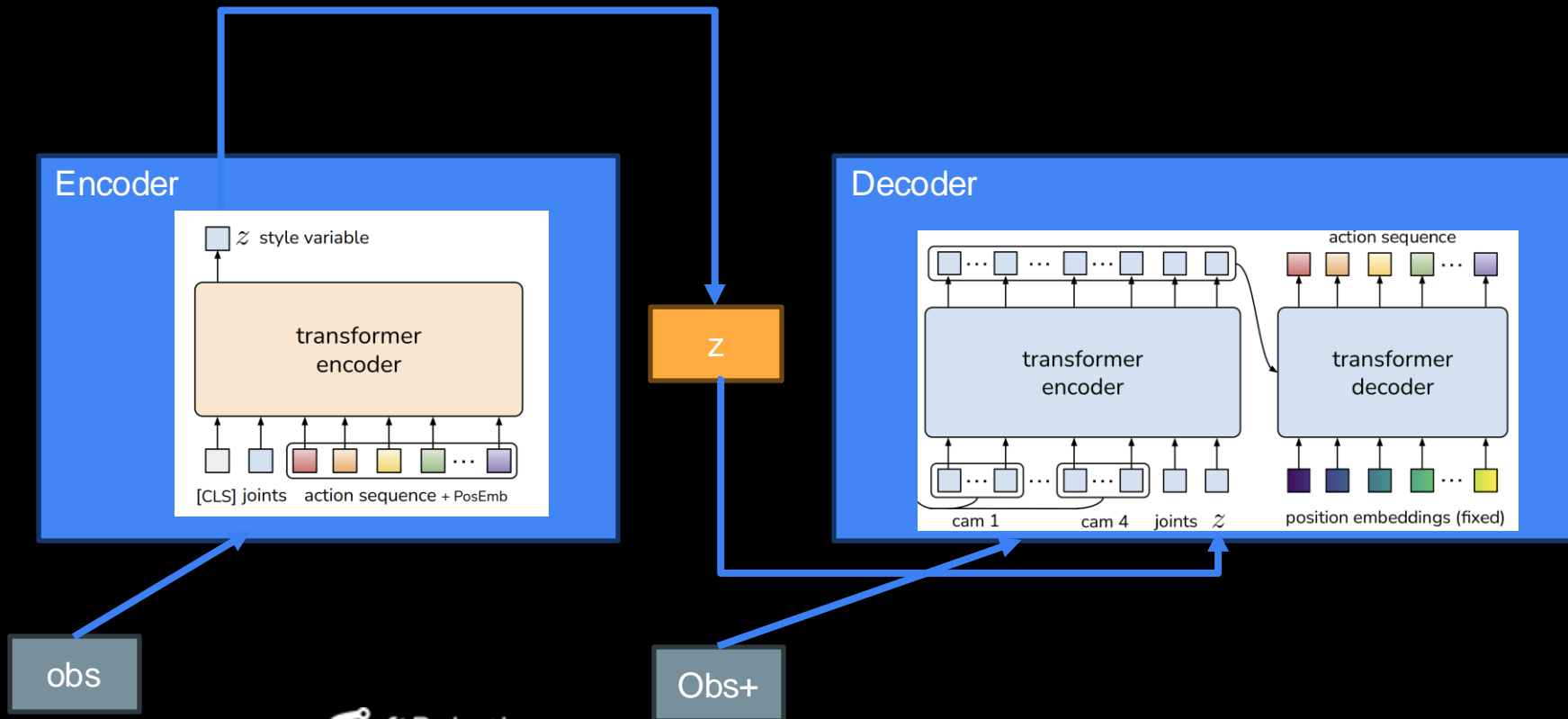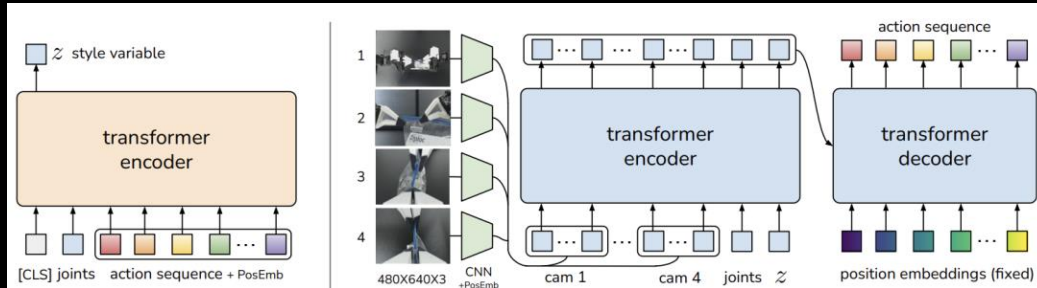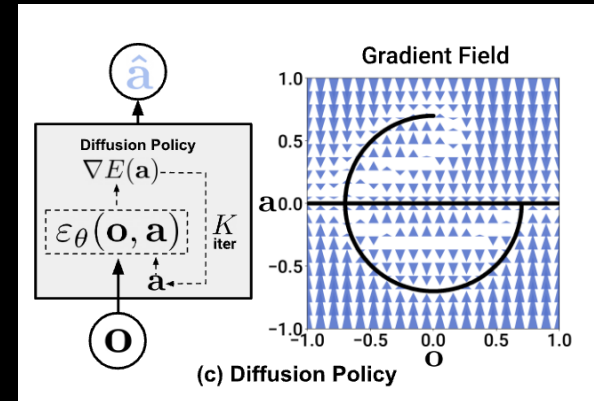
# Action Chunking Transformers

# Contents

What we need to understand today:
1. Transformer Architecture
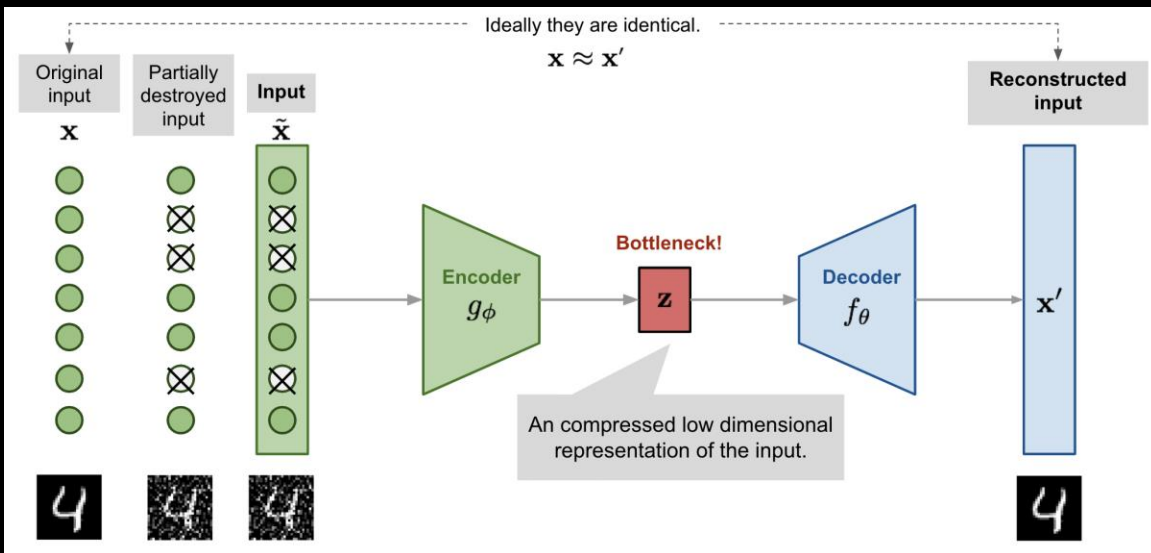2. Variational Autoencoder
3. Diffusion Mechanism



ACT



Diffusion Policy

# Denoising

Besides compressing, vae models can be trained to recover the original input.
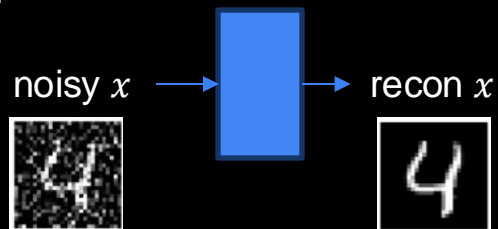


Resource: https://lilianweng.github.io/posts/2018-08-12-vae/

```python
def train_iter(self, x):
    x_noise = x + torch.rand_like(x)
    recon_x = self.forward(x)
    loss = F.mse_loss(x, recon_x)
    loss.backward()
```
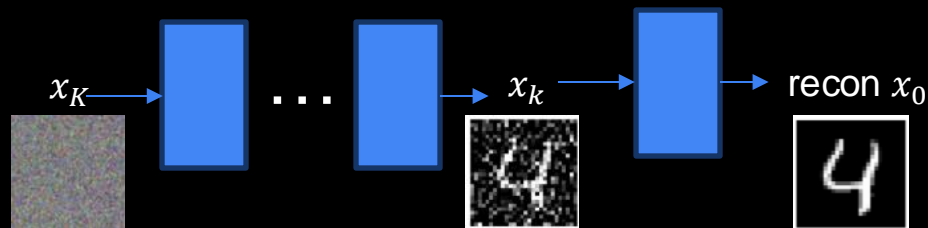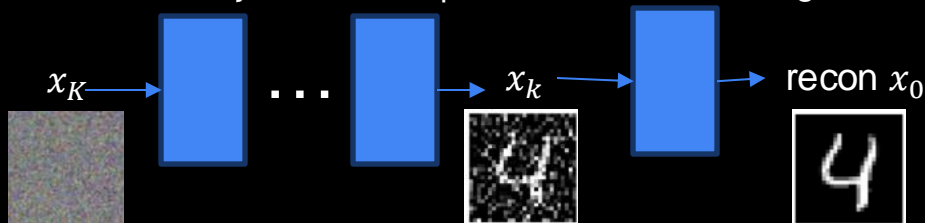
# Diffusion Models

- Given a denoise model

noisy $x$ → recon $x$

- What about stacking them many times

$x_K$ → ... → $x_k$ → recon $x_0$

ETH zürich    SoftRobotics Laboratory

# Diffusion Models

- We stack many denoise steps to reconstruct from gaussian noise

$x_K \longrightarrow$ [ ] $\ldots$ [ ] $\longrightarrow x_k \longrightarrow$ [ ] $\longrightarrow$ recon $x_0$

- How to train the network
  - aka: how to define the loss          Loss=MSE($x_k$  $\widetilde{x_k}$)
    - aka: how to find $x_k$ and $x_{k+1}$
- $x_{k+1}$ is can be computed by adding noise from $x_k$ (**Forward diffusion process**)

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

- But we only have $x_0$
  - Good news: $x_{k+1}$ can be directly computed from $x_0$

Hugging Face
https://huggingface.co › api › schedulers › overview ⋮
**Schedulers**
**Diffusers** provides many **scheduler** functions for the diffusion process. A **scheduler** takes a model's output (the sample which the diffusion process is iterating ...

# Diffusion Models: Train with Scheduler

```python
noise_scheduler = DDPMScheduler(num_train_timesteps=1000)
def train_iter(batch):
    clean_images = batch["images"]
    # Sample noise to add to the images
    noise = torch.randn(clean_images.shape).to(clean_images.device)
    bs = clean_images.shape[0]
    # Sample a random timestep for each image
    timesteps = torch.randint(
        0, noise_scheduler.config.num_train_timesteps, (bs,),
    ).long()
    # Add noise to the clean images according to the noise magnitude at each timestep
    # (this is the forward diffusion process)
    noisy_images = noise_scheduler.add_noise(clean_images, noise, timesteps)
    # Predict the noise residual
    noise_pred = model(noisy_images, timesteps)
    loss = F.mse_loss(noise_pred, noise)
```

# Diffusion Models: Generate with Scheduler

```python
# Initialize random noise as the starting point for generation
target = torch.randn((num_samples,) + image_shape)
# Iteratively remove noise step by step
for t in noise_scheduler.timesteps:
    # Model prediction for noise residuals
    with torch.no_grad():
        noise_pred = model(target, t).long())
    # Compute the previous image state by reversing the diffusion process
    target = noise_scheduler.step(noise_pred, t, target).prev_sample
```