





17 Nov

24 Nov

1 Dec

Imitation Learning

14:15-15:00

Davide Liconti

Guest Lecture

*Oier Mees @
Microsoft Spatial AI Lab*

Product Development

*Benedek Forrai @
Mimic Robotics*

Training Policies with SRL_IL (workshop)

15:15-16:00

Chenyu Yang

VLAs and Foundation Models

Davide Liconti

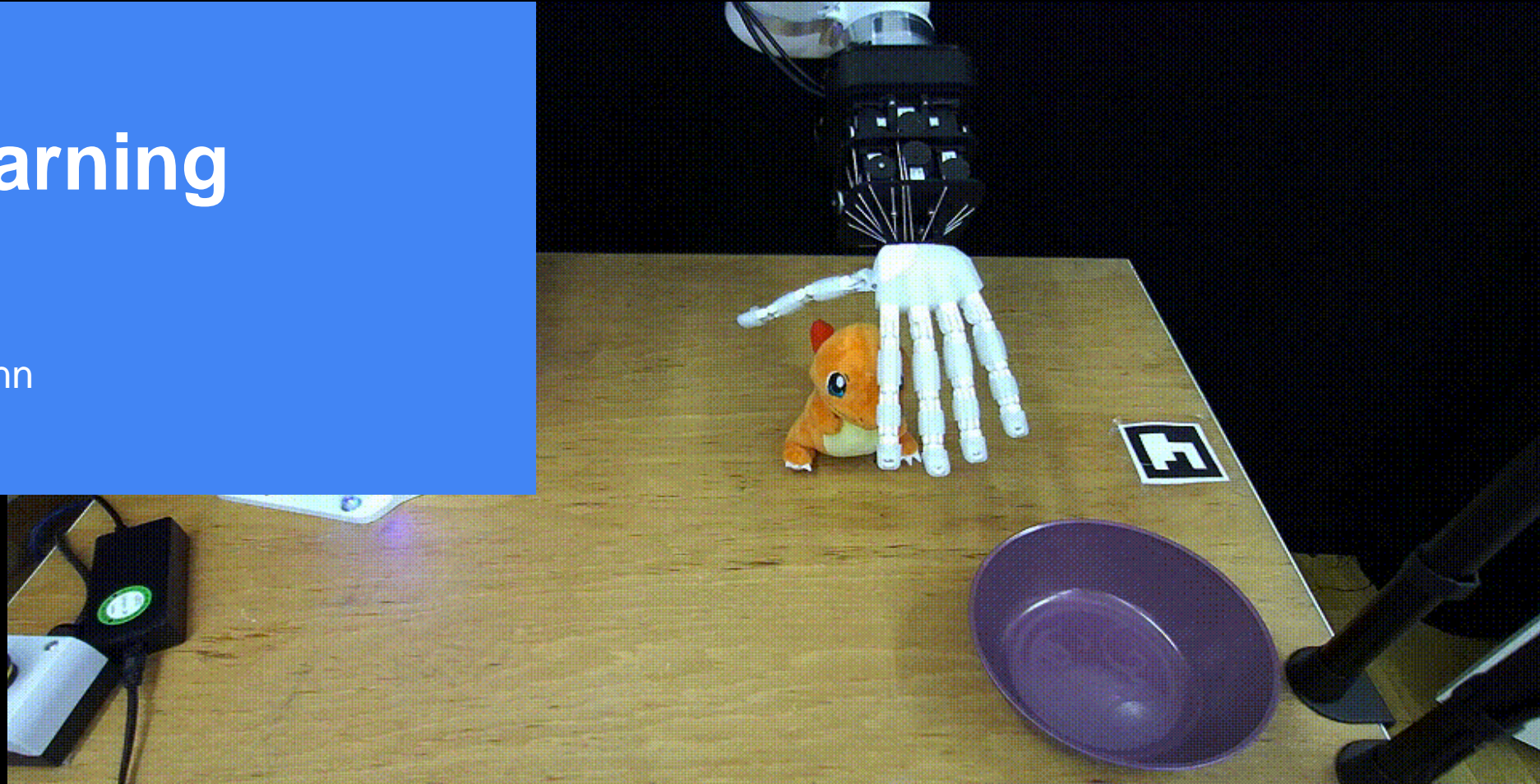
Class Wrap-up and Q&A

Robert Katzschmann



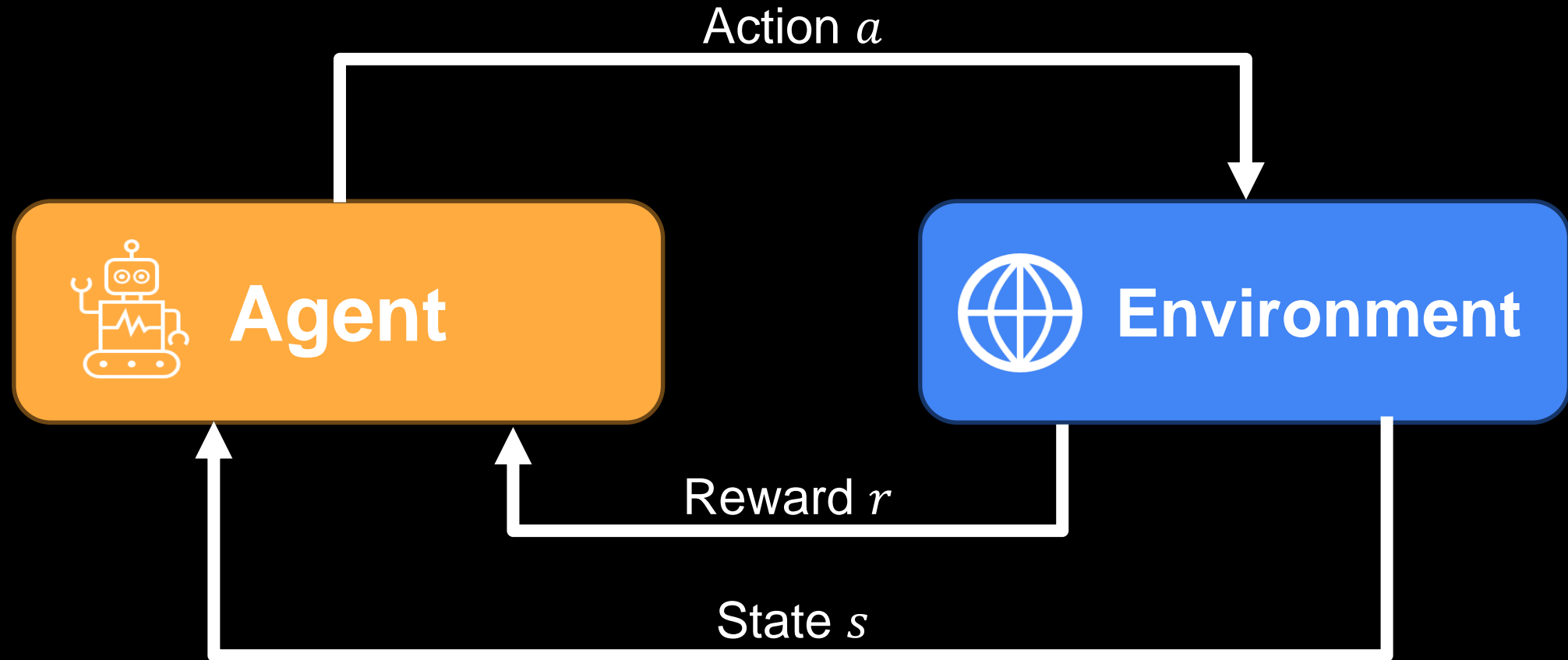
Imitation learning

Davide Liconti
Prof. Robert Katzschmann
17 November 2025





Reinforcement Learning (RL) Recap



How to define the reward?

What's the reward function for this task?

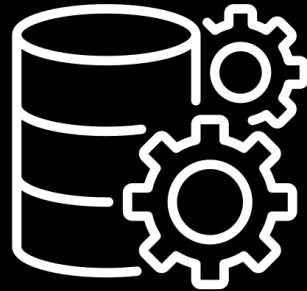


- **Sparse reward:** +1 when tied (RL won't explore enough).
- **Dense reward:** but what is the distance to “tiedness”?
- Multi-step, long-horizon manipulation → reward shaping quickly becomes *arbitrary* or *impossible*.



Imitation Learning (IL)

aka Behavior Cloning (BC), aka Learning from Demonstrations (LfD)



Given a dataset of “Expert transitions”

$$\mathcal{D} = \{(o_t, a_t)\}_{t=1}^N$$

Simple supervised regression : $\pi = \arg \min \mathbb{E}_{(o,a) \sim \mathcal{D}} [\|a - \pi(o)\|^2]$

How to collect data? → previous lecture



Imitation Learning (IL)

why does it dominate today

Today **IL** is by far the dominant approach for manipulation

- **No** complex or arbitrary **reward shaping**
- **No sim2real gap** (for real world teleop data)
- Simpler, easier to debug and interpret
- Can theoretically **scale** with data and compute

LLMs (e.g., GPT) are trained in a similar form as IL.
They are trained to predict the next token given a "context" of recent tokens.

Why is learning actions different than learning next token?



Imitation Learning (IL)

end-to-end learning

Traditional Perception + Control



End-to-end Policies



Visuomotor Policies



End-to-end Policies



Types of Observations

RGB Images



Others

- *Depth*
- *Tactile*
- ***Proprioception***
- *etc.*

Issues with Observations

Pixel-space is a terrible space to do control

- High dimensional
- Discrete and non-differentiable

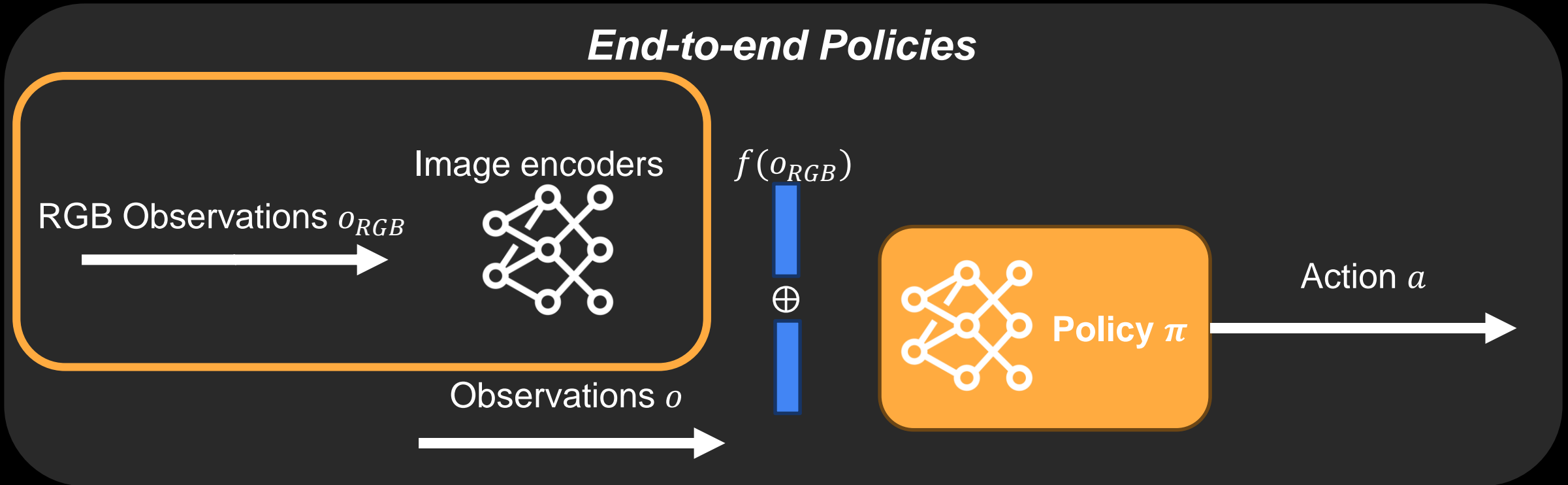


Need latent, meaningful features

Visuomotor Policies: Encoding Observations



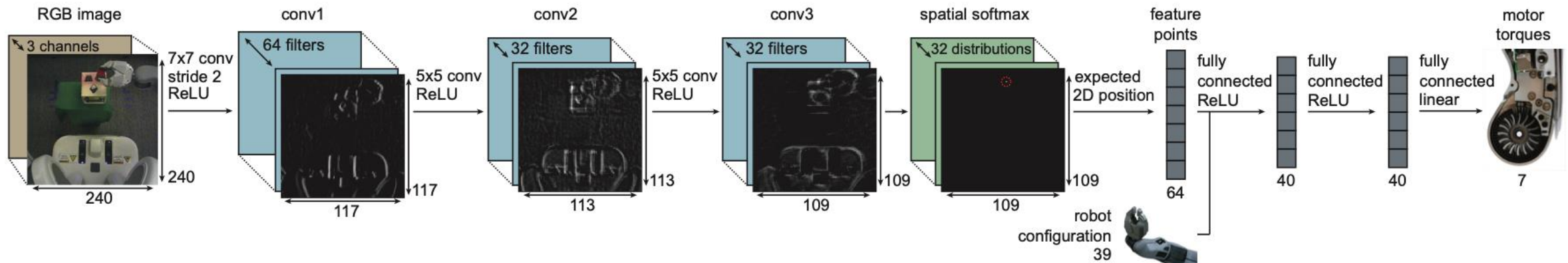
End-to-end Policies



Deep Visuomotor Policies

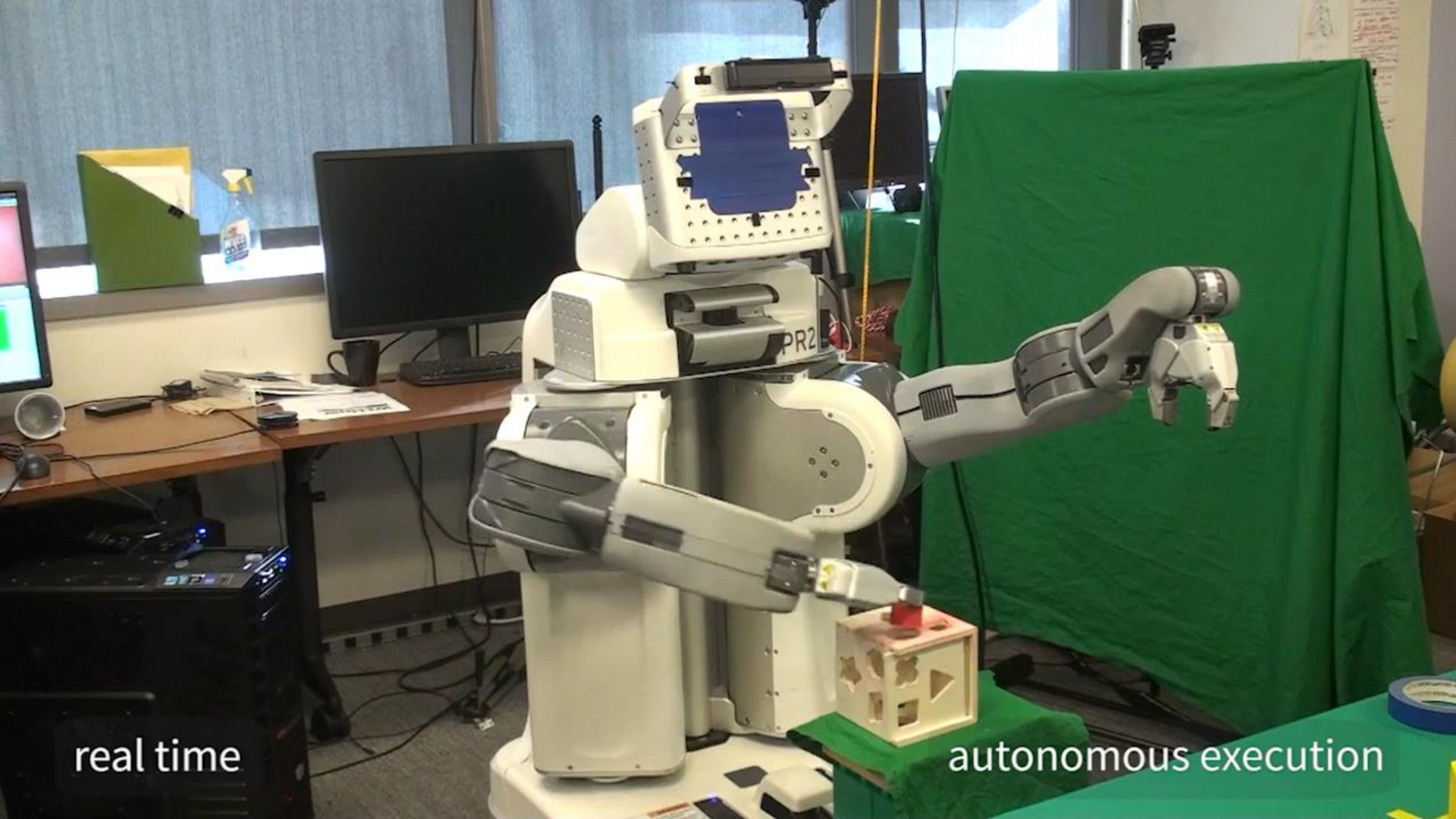


Sergey Levine*, Chelsea Finn*, Trevor Darrell, Pieter Abbeel, "End-to-end training of deep visuomotor policies", *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 1334--1373, 2016.



The name "visuomotor" is chosen to emphasize that the policy is trained to predict actions directly from RGB camera images

The action space was directly motor torques. Today, higher level actions are used, and let IK + low-level control transfer to torques



real time

autonomous execution

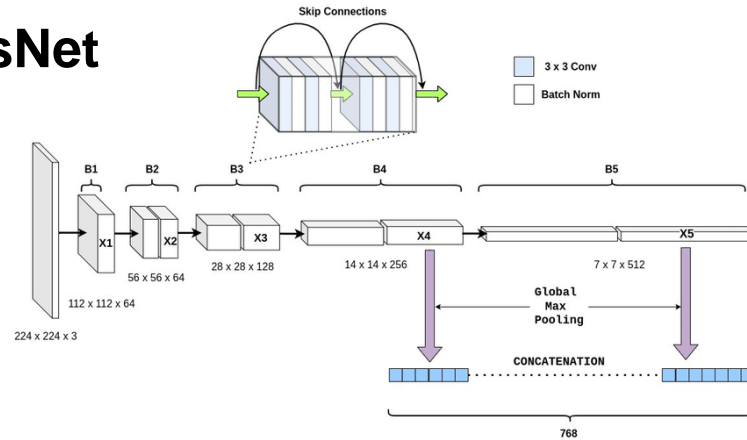
Visual Encoders Trained on Large Data

IMAGENET

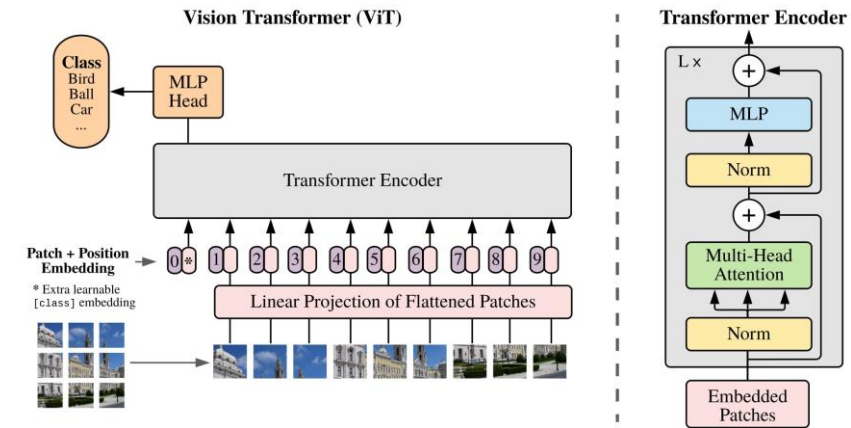


Today, a lot of visual encoders trained on large image datasets

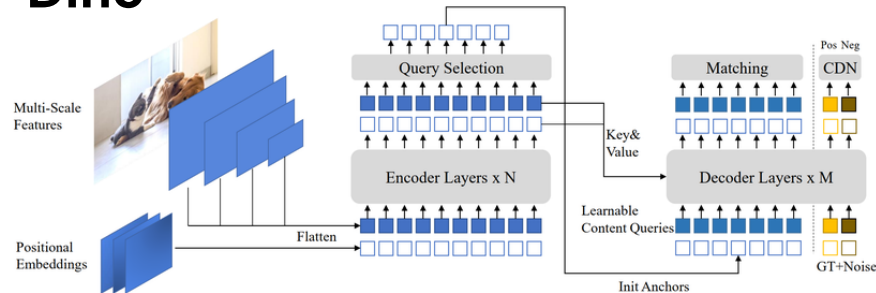
ResNet



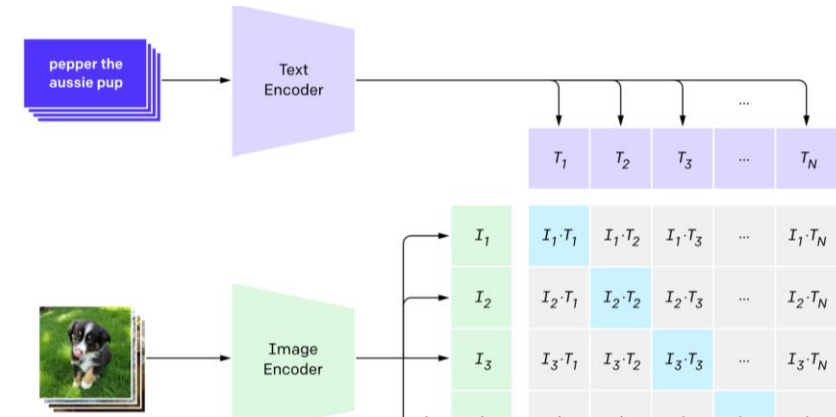
ViT



Dino

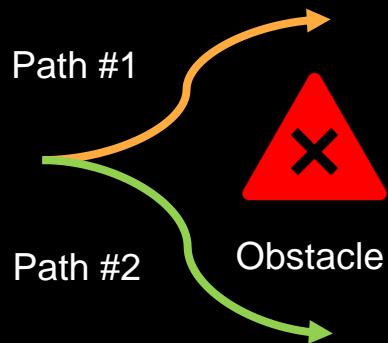
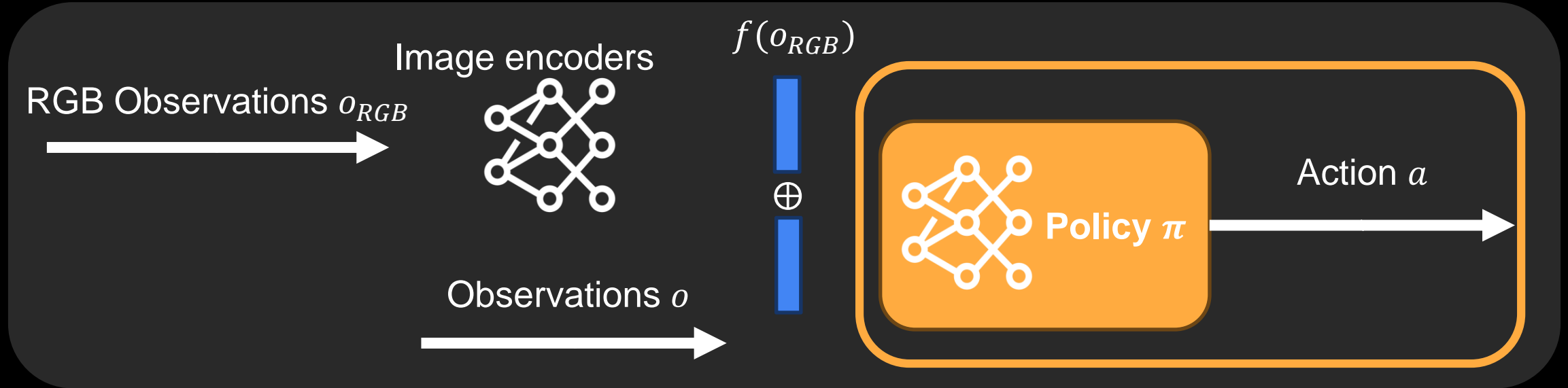


CLIP



We assume $f(o_{RGB,t}) = \text{image_encoder } o_{RGB,t}$.
The policy maps $\pi: f(o_{RGB}) + o \rightarrow a$

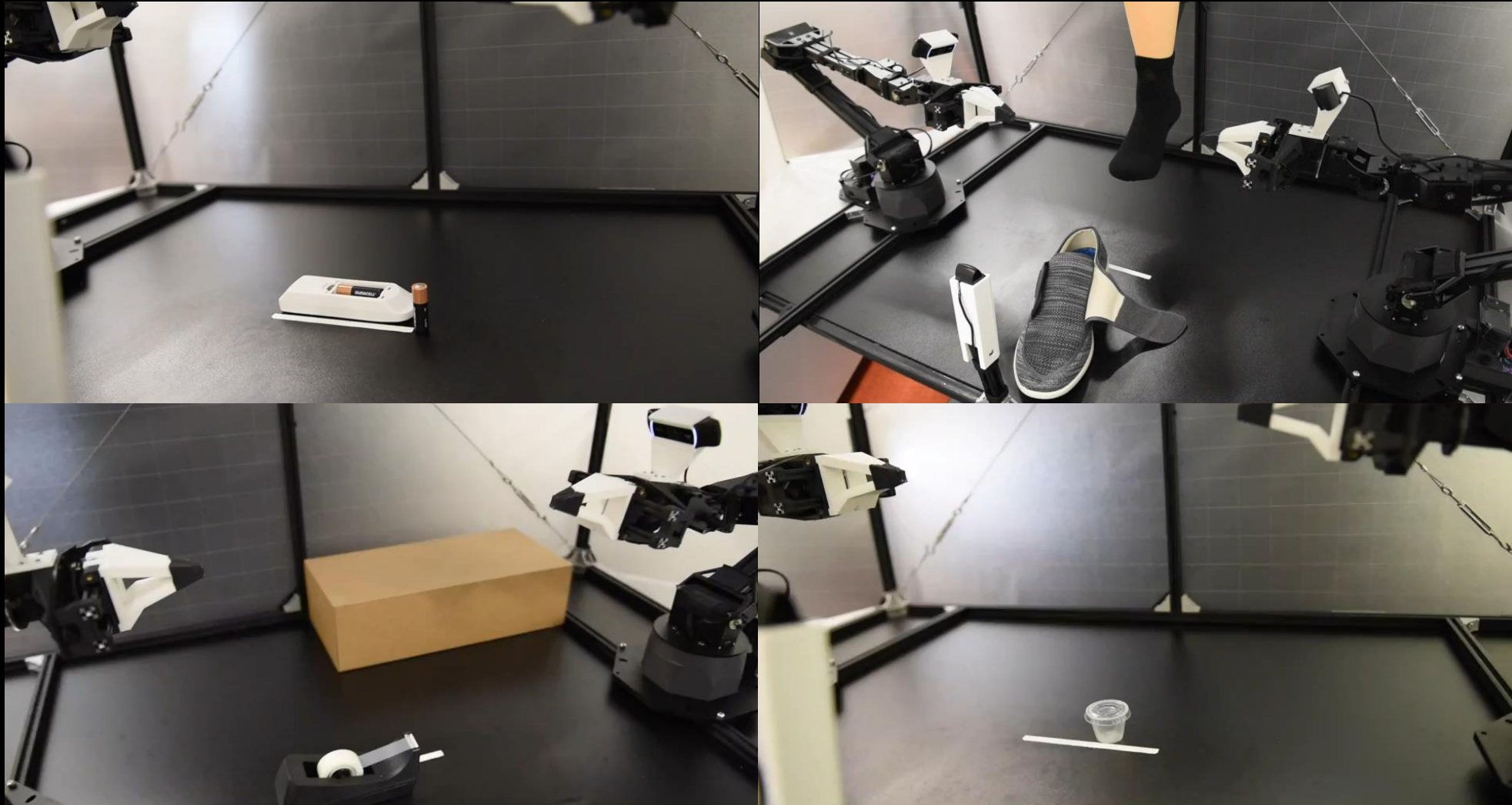
Visuomotor Policies : Predicting Actions



What is the best way to predict actions?

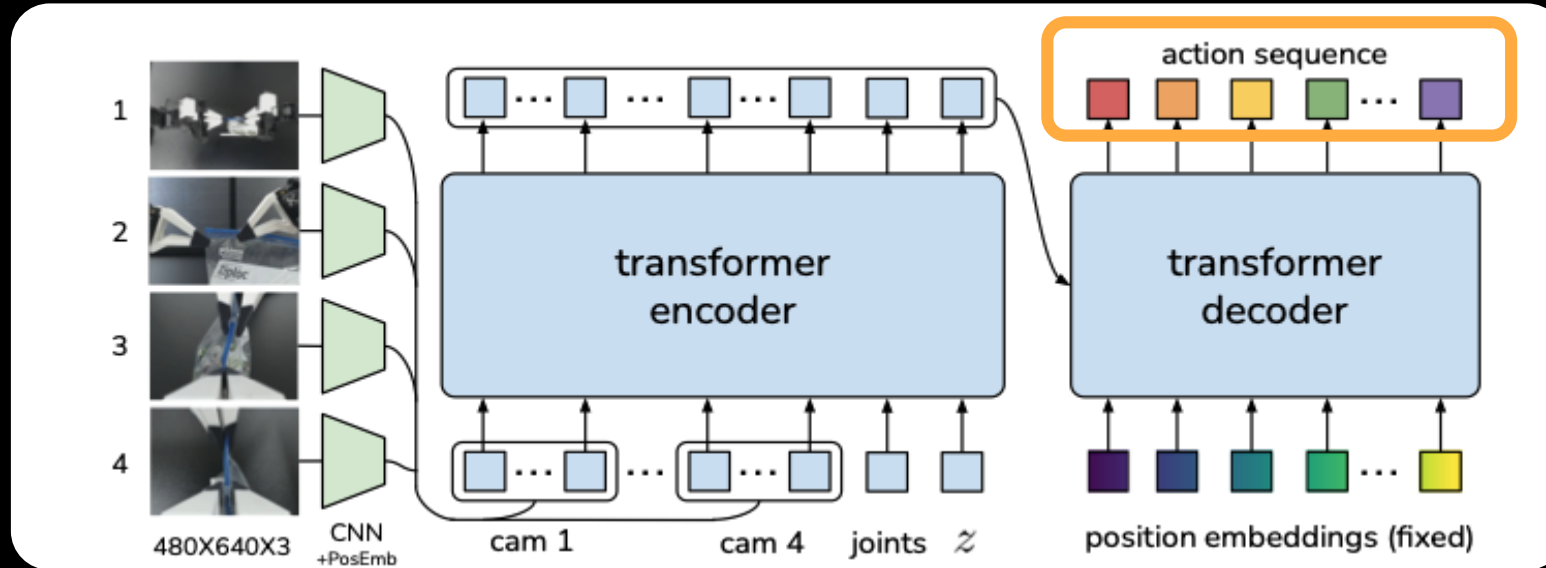
- Multimodality: for the same observation there can be multiple valid actions
- In 2023, **ACT** (Action Chunking Transformers) and **Diffusion Policy** showed convincing capabilities to learn complex dexterous actions

ACT (Action Chunking Transformers)





ACT (Action Chunking Transformers)



Predict a **chunk** of k actions, instead of a single step

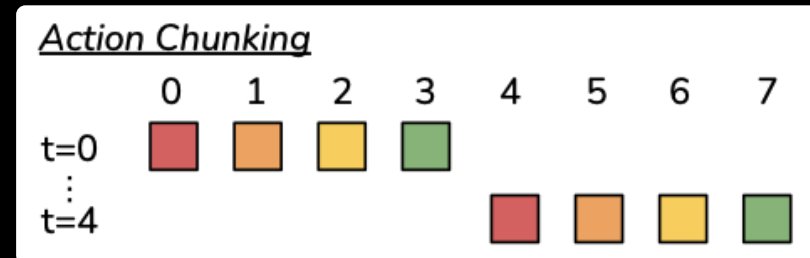
$$\pi_{\theta}(a_t|o_t) \rightarrow \pi_{\theta}(a_{t:t+k}|o_t)$$

- Avoid compounding error (wrong predictions quickly leads to out of distribution)
- Shorter horizon control
- Help model non-Markovian behavior (e.g., pauses in the middle of a demonstration)



ACT : Temporal Ensemble

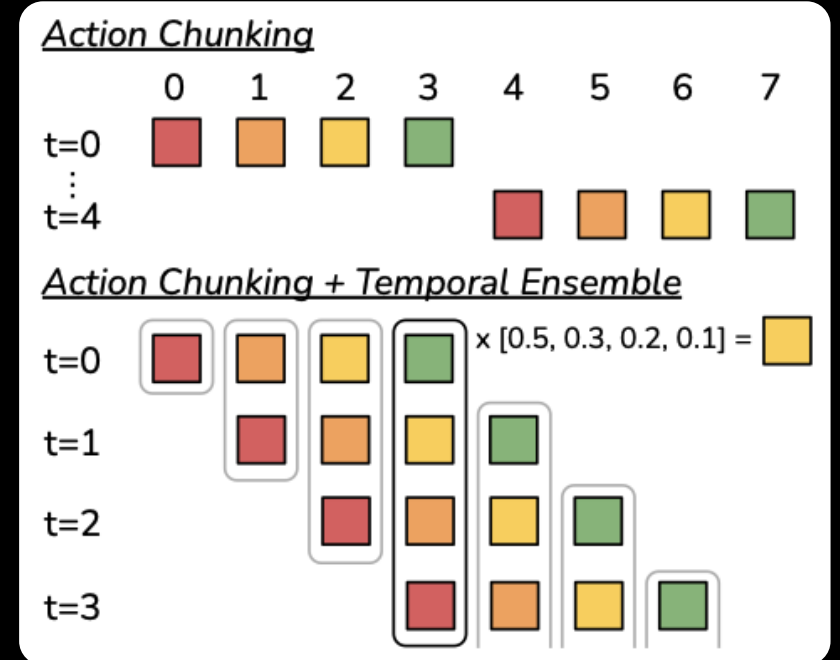
- Naïve implementation of action chunking can be suboptimal
- If new observation is incorporated abruptly every k steps this can result in jerky robot motion





ACT : Temporal Ensemble

- Naïve implementation of action chunking can be suboptimal
- If new observation is incorporated abruptly every k steps this can result in jerky robot motion



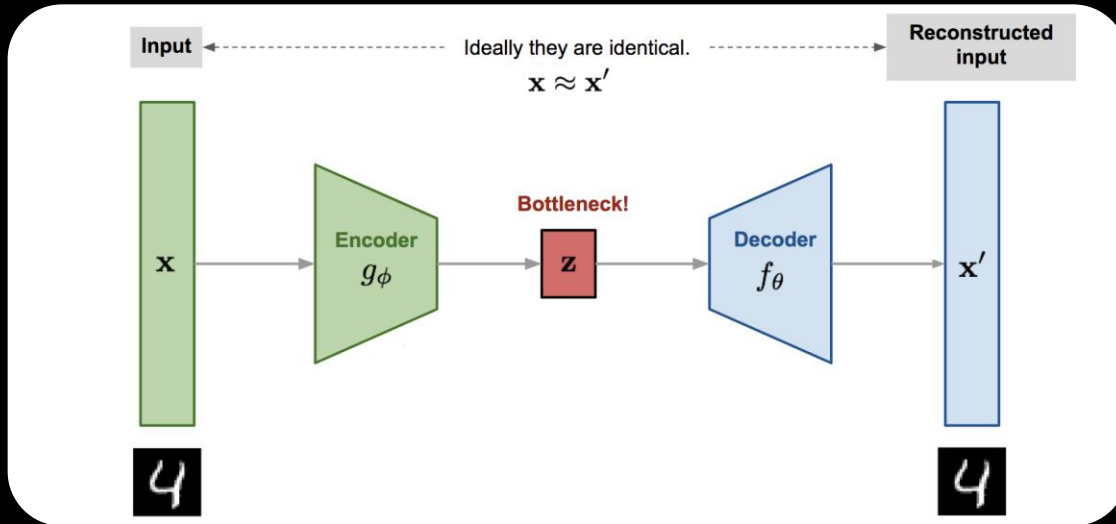
Temporal ensemble

Query the policy at every timestep, aggregate them with a weighted exponential average

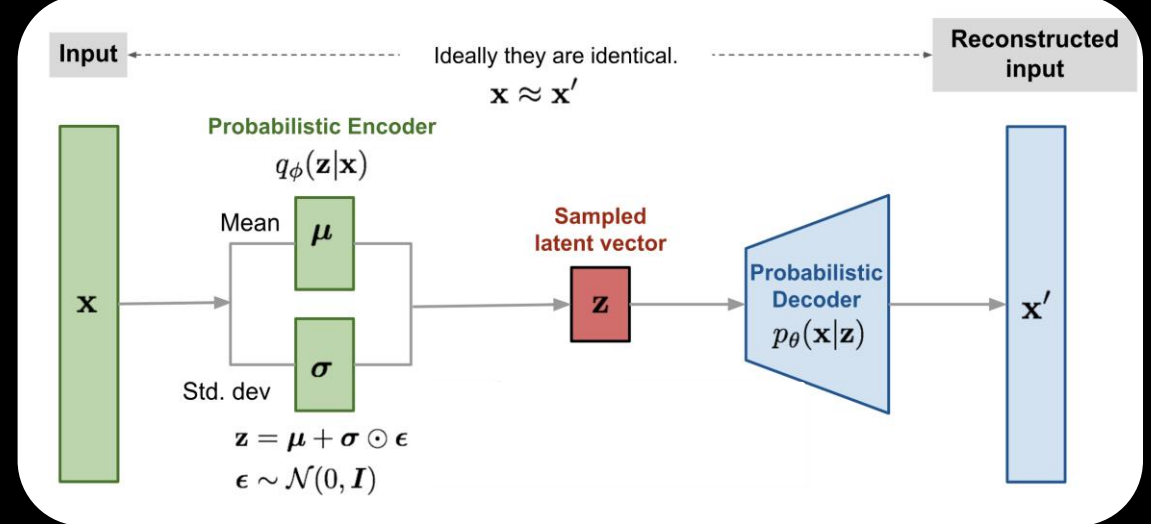
$$w_i = \exp(-m * i)$$

ACT : VAEs

AutoEncoder

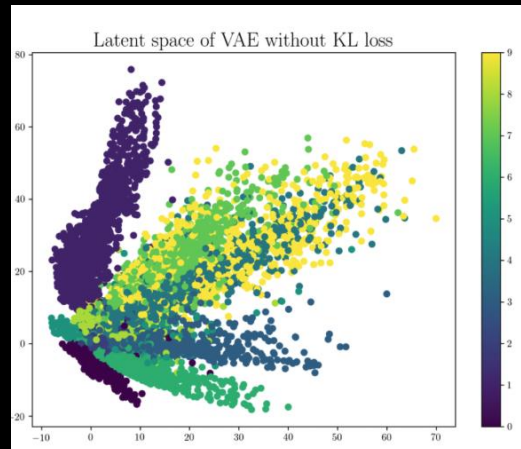


Variational Autoencoder

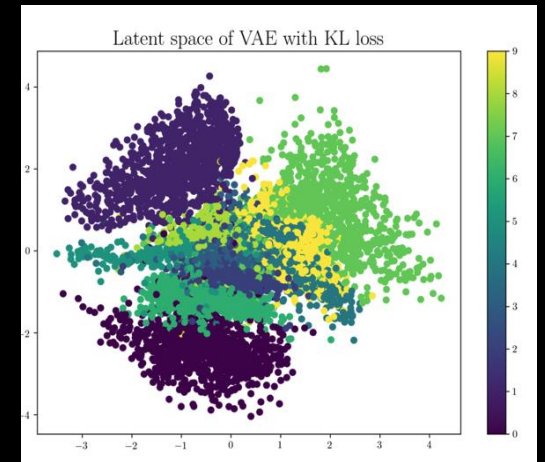


$$L = \|f(g(x)) - x\|$$

$$= \|x' - x\|$$



$$L = \underbrace{\|x' - x\|}_{\text{reconstruction}} + \underbrace{D_{KL}(N(\mu_x, \sigma_x) | N(0, 1))}_{\text{regularization}}$$

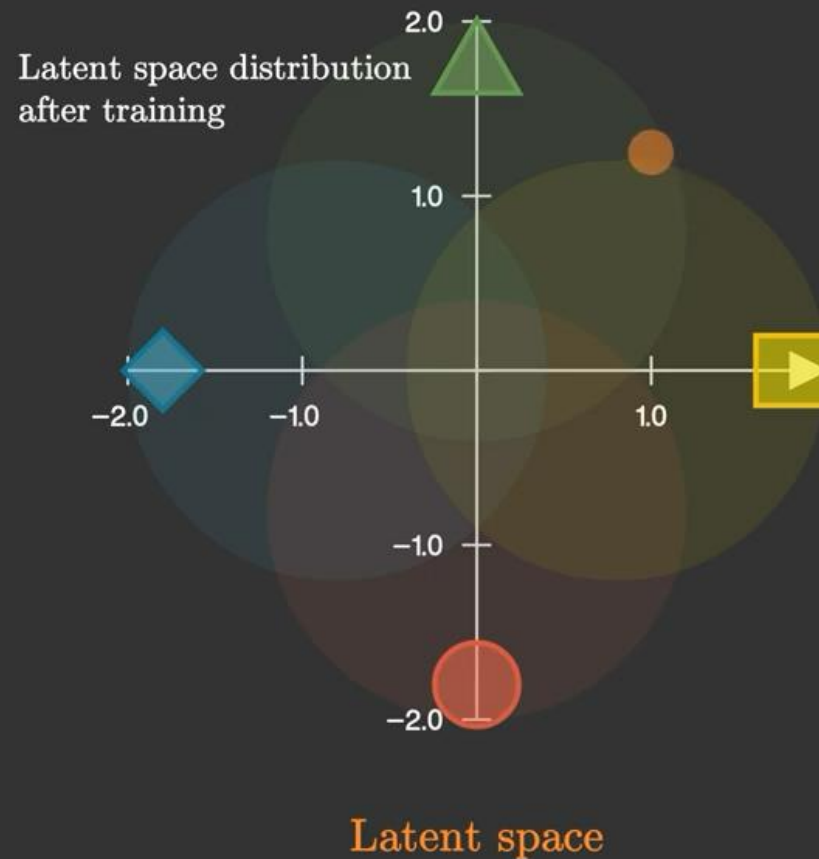


ACT : VAEs



Latent space of Variational Autoencoder

aqeel



Latent space is regularized. Vectors sampled from latent space can generate valid data.



Vectors sampled from overlapping distribution generates morphed data.

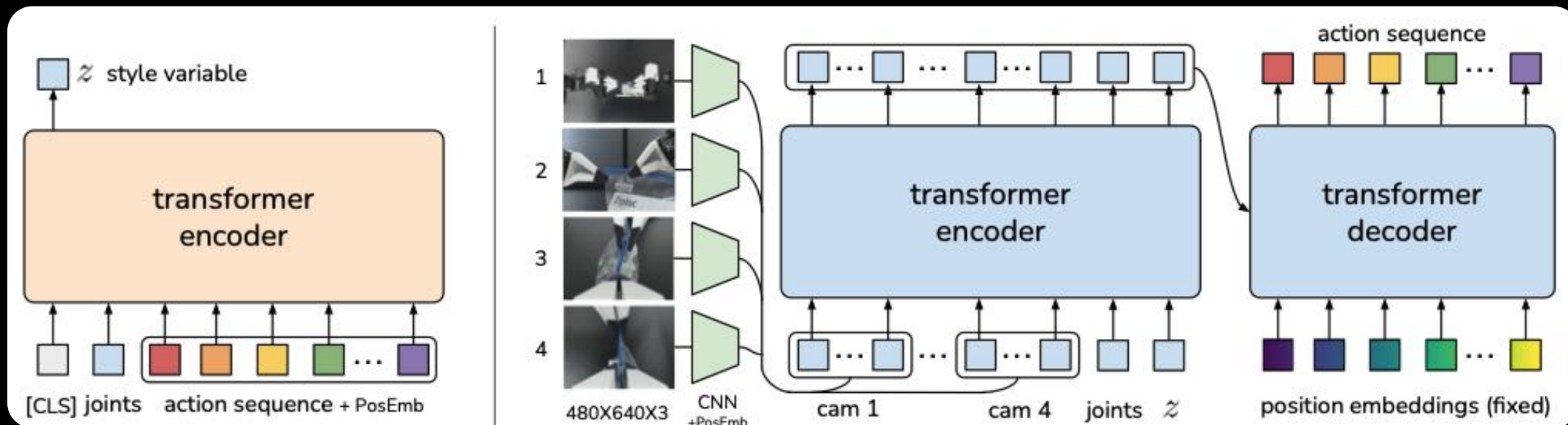
Decoder's output

More on: www.MLinGIFS.aqeel-anwar.com

ACT: CVAE



Train the policy as a conditional variational autoencoder (CVAE)



Encoder

$$q_{\phi}(z | a_{t:t+k}, \bar{o}_t)$$

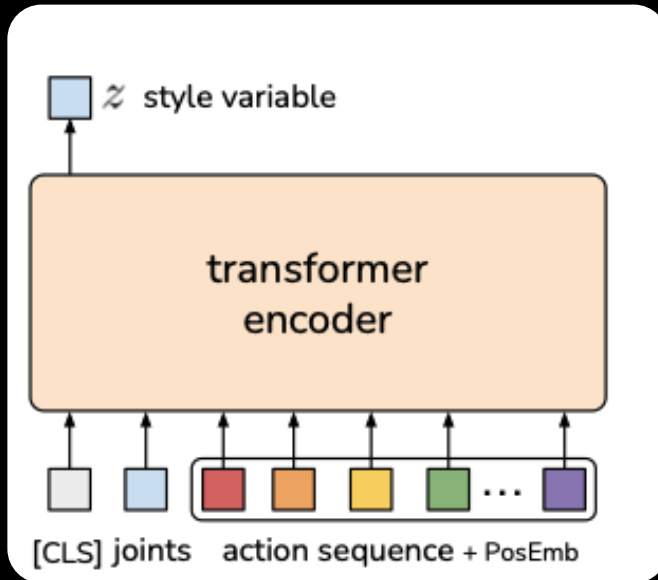
Conditional Decoder

$$\pi_{\theta}(a_{t:t+k-1} | o_t, z)$$



ACT : CVAE

Train the policy as a conditional variational autoencoder



At training time, we encode the action sequence and the proprioception using a transformer

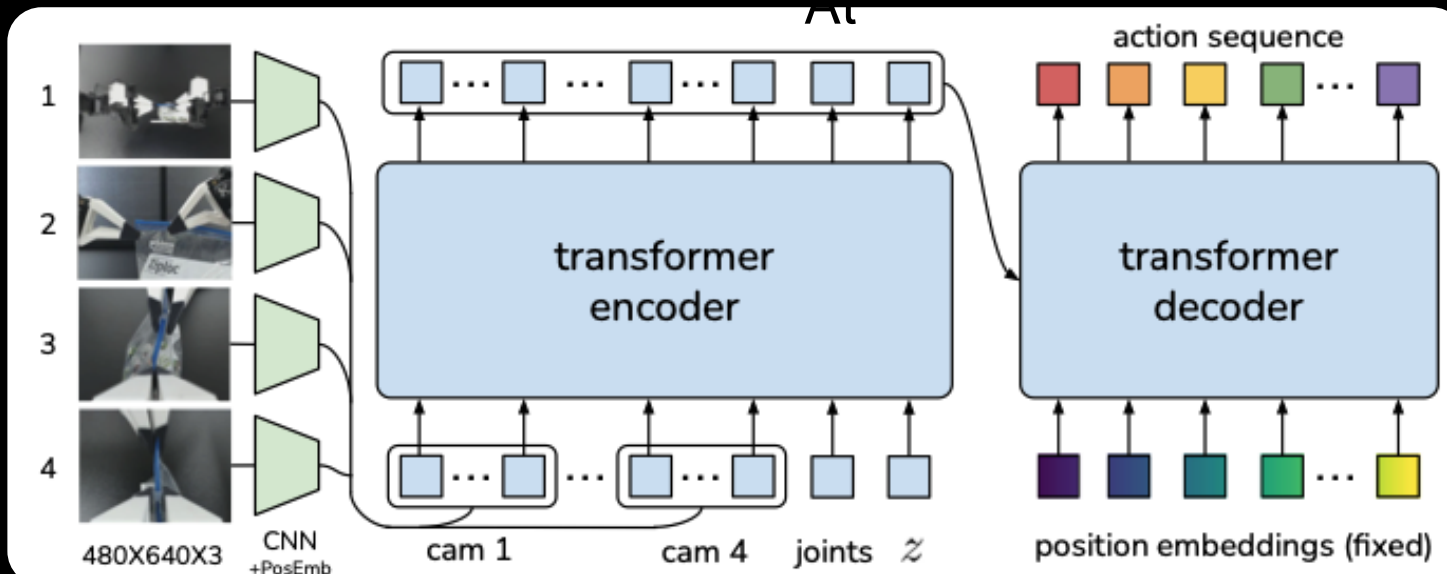
The output is the encoded representation $\mathbf{z} = q_{\phi}(z | a_{t:t+k}, \bar{o}_t)$

$$\bar{o}_t = o_t - o_{RGB}$$



ACT: CVAE

Train the policy as a conditional variational autoencoder



Then, we decode the encoded vector z conditioned on image features and proprioception.

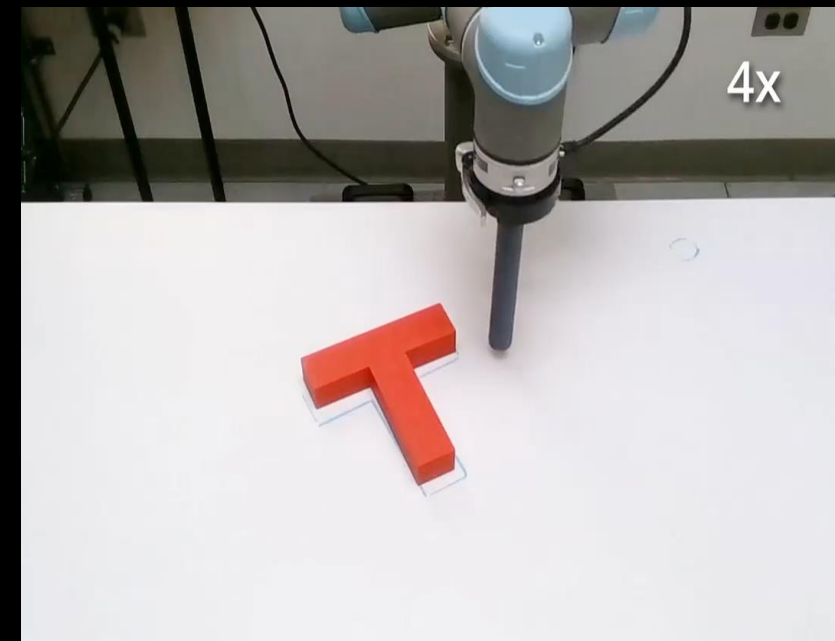
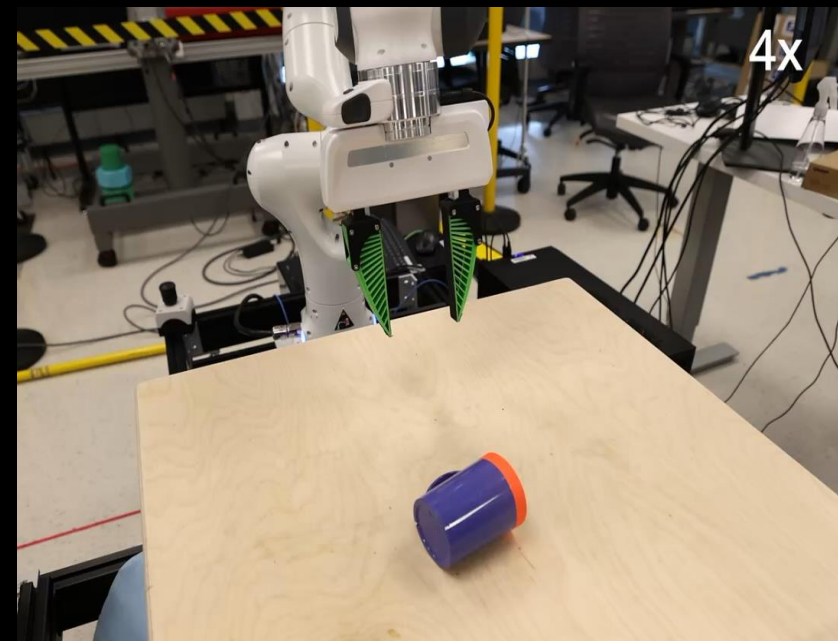
$$a_{t:t+k} = \pi_{\theta}(a_{t:t+k-1} | o_t, z)$$

Training objective is the VAE loss

$$\|a_{t:t+k} - \hat{a}_{t:t+k}\|^2 + \beta * D_{KL}(\mathcal{N}(\mu_z, \sigma_z), \mathcal{N}(0, 1))$$

At inference time z is set to 0 (mean of the distribution), since we don't have access to future actions

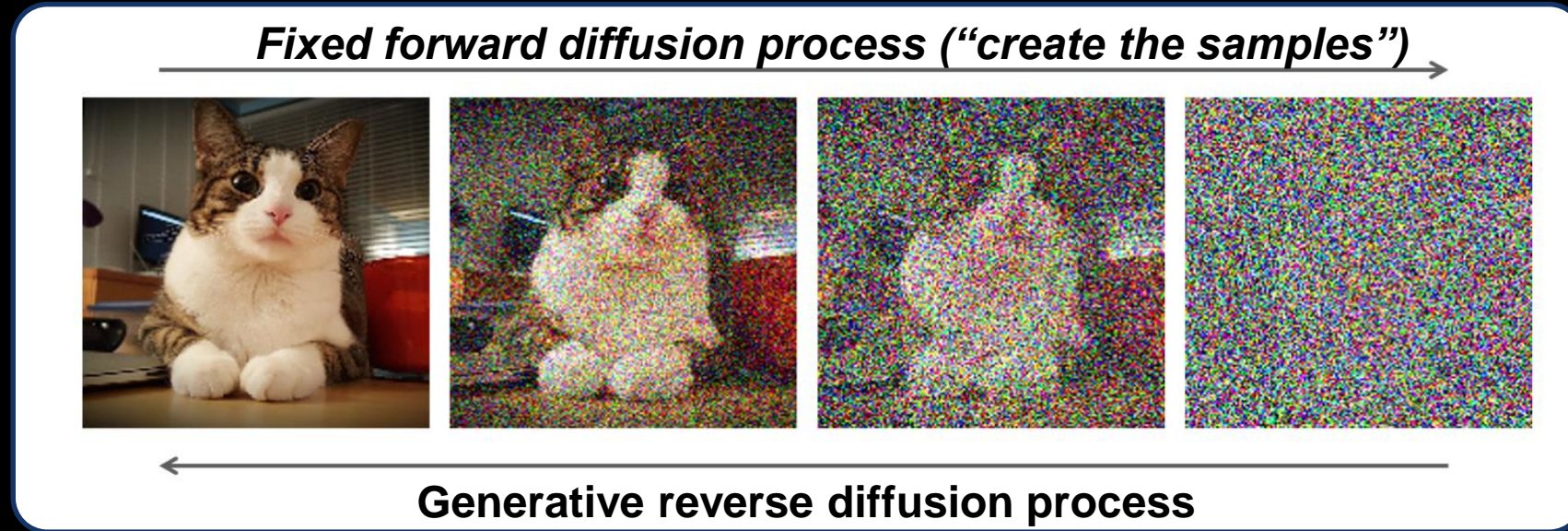
Diffusion Policy



Chi, C. et al. (2023). *Diffusion Policy: Visuomotor Policy Learning via Action Diffusion*



Diffusion Policy – Diffusion Process



1. Sampling $x_0 \sim \mathcal{D}$, noise level $\sigma \sim [\sigma_{min}, \sigma_{max}]$, noise $\epsilon \sim N(0, I)$
2. Generating noisy data $x^\sigma = x^0 + \sigma \epsilon$
3. Predicting ϵ (direction of noise) from x^σ by minimizing squared loss

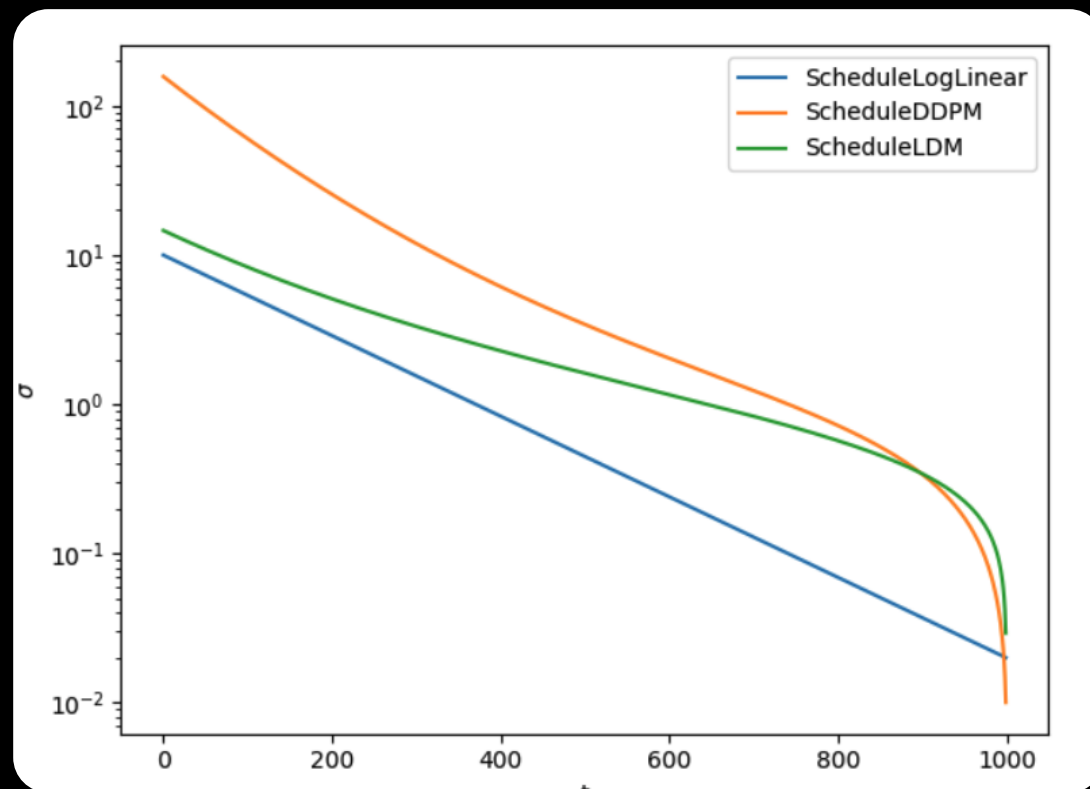
This amounts to training a θ -parametrized neural network $\epsilon_\theta(x, \sigma)$, by minimizing the loss function $L(\theta) = \mathbb{E} \|\epsilon_\theta(x^0 + \sigma_t \epsilon, \sigma_t) - \epsilon\|^2$



Diffusion Policy – Diffusion Process

How do we add noise?

- In practice, σ is not sampled uniformly in the interval $[\sigma_{min}, \sigma_{max}]$, but the interval is discretized into K values called **schedules** and then sampled uniformly between the schedules.
- If we now follow the noise schedule, we can find the denoising formulation
- A well-known scheduler is the **ScheduleDDPM** (Denoising Diffusion Probabilistic Models)



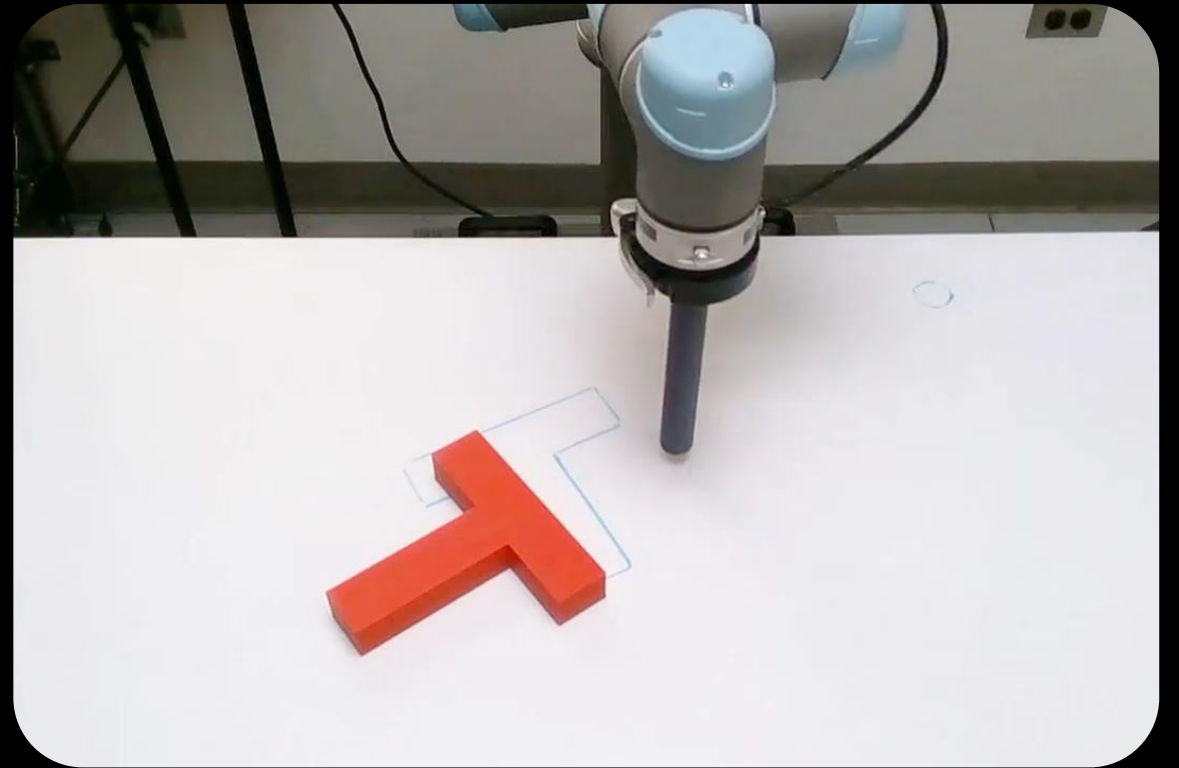


Diffusion Policy

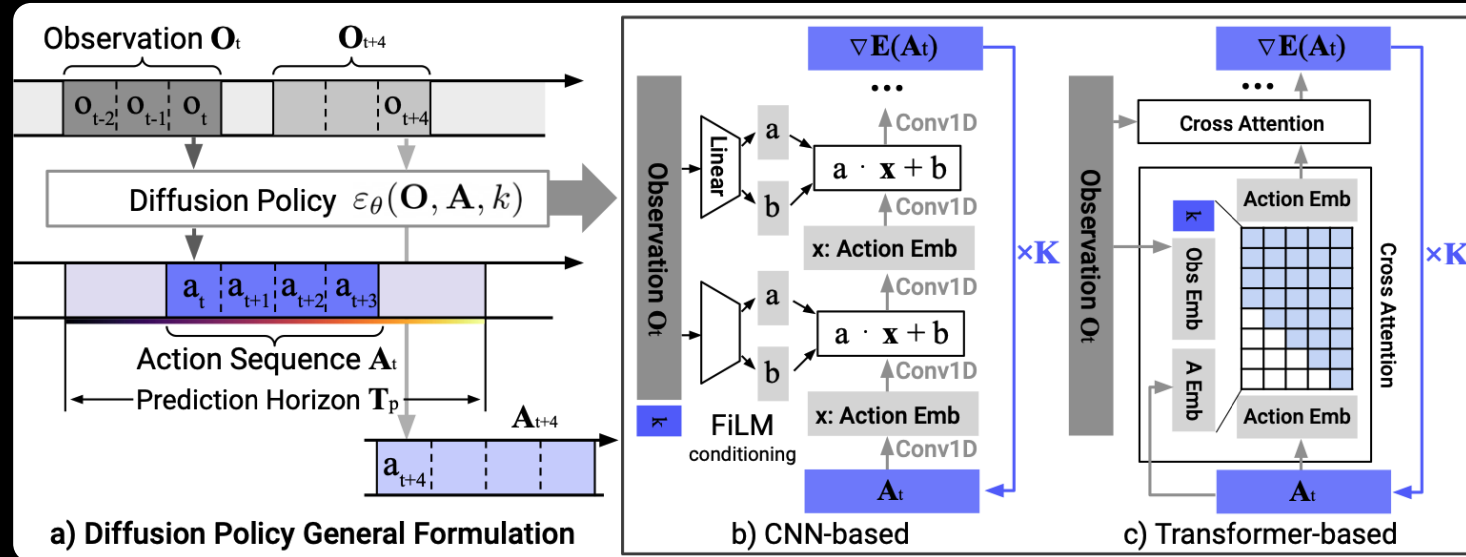
Apply diffusion process to robot action prediction

Advantages

- Handles multimodal action distribution (many valid actions for a given observation)
- Works in high-dimensional space (e.g., image generation)
- Stable training



Diffusion Policy



Training

$$L(\theta) = \|\epsilon^k - \epsilon_\theta(O_t, A_t^0 + \epsilon^k, k)\|^2$$

O_t : **history** of observations

A_t^0 : future action sequence

k : noise timestep (scheduler)

t : action timestep

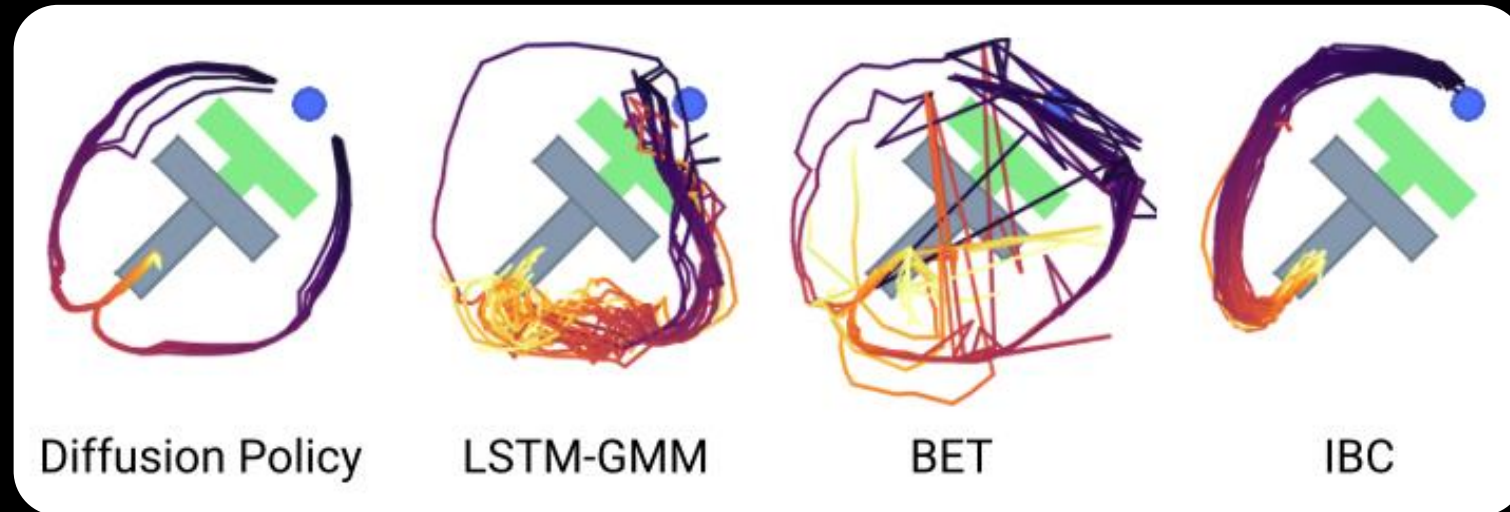
Denoising Process

We use a DDPM to approximate the conditional distribution $p(A_t | O_t)$

$$A_t^{k-1} = \alpha \left(A_t^k - \gamma \epsilon_\theta(O_t, A_t^k, k) + \mathcal{N}(0, \sigma^2 I) \right)$$



Diffusion Policy - Multimodality



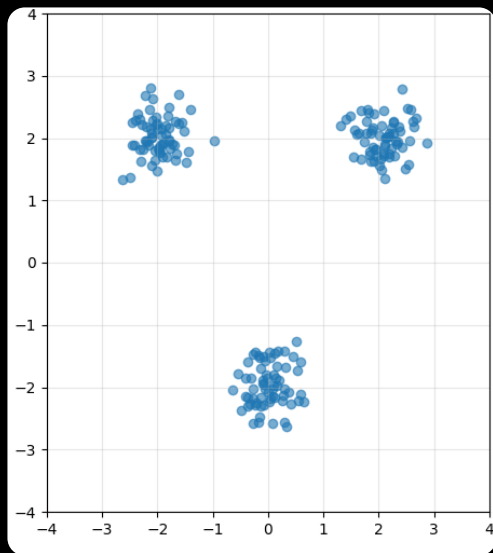
Diffusion Policy naturally models multimodal action distributions:

- Each rollout **commits to a coherent action mode**, but different rollouts represent **different valid strategies** - capturing true multimodality.

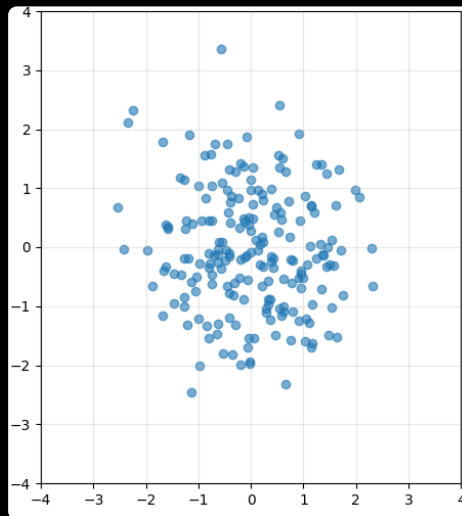


Side quest: Flow Matching

Real data distribution

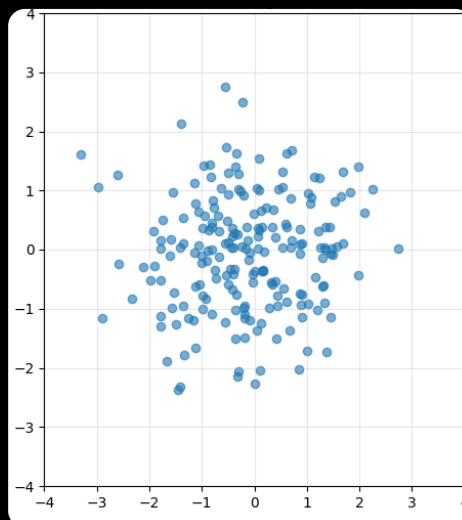


Diffusion



Diffusion Models gradually add noise to data until it becomes pure noise, then learn to reverse this process

Flow Matching



Flow matching creates a continuous path (or flow) between noise and data distribution



Side quest: Flow Matching

```
# DDPM: learn to predict noise
for each training step:
```

```
    x0 = sample_data()
    t = random_timestep()
```

```
    # Add noise
```

```
    eps = random_noise()
```

Noise scheduler

```
    xt = add_noise(x0, eps, t)
```

```
    # Predict the noise
```

```
    eps_pred = model(xt, t)
```

```
    # Train to match the true noise
```

```
    loss = MSE(eps_pred, eps)
```

```
    update(model, loss)
```

```
# Flow Matching: learn velocity from noise → data
for each training step:
```

```
    x0 = sample_data()
```

```
    z = random_noise()
```

```
    t = random_time_between_0_and_1()
```

```
    # Point along straight line from noise to data
```

```
    xt = (1 - t) * z + t * x0
```

```
    # Target velocity field
```

```
    v_target = x0 - z
```

```
    # Predict velocity
```

```
    v_pred = model(xt, t)
```

```
    # Train to match the velocity
```

```
    loss = MSE(v_pred, v_target)
```

```
    update(model, loss)
```



Key Takeaways

- Imitation Learning turns control into supervised learning
- Pixels become useful through feature encoders
- Demonstrations are noisy & multimodal
→ generative models (CVAE in ACT, diffusion in DP) capture variability and improve precision in challenging manipulation tasks.

Next lecture : From single-task to multi-task? What happens if we scale the model up?

VLA models



Acknowledgements and useful resources

- MIT 6.8210 Underactuated Robotics, chapter 21 <https://underactuated.mit.edu/imitation.html>
- Diffusion models :
 - (easy): <https://chenyang.co/diffusion.html>
 - (more complex) <https://lilianweng.github.io/posts/2021-07-11-diffusion-models/>
- Flow matching vs Diffusion <https://harshm121.medium.com/flow-matching-vs-diffusion-79578a16c510>
- Diffusion and Flow models lecture <https://arxiv.org/pdf/2506.02070>
- Levine, S., Finn, C., Darrell, T., & Abbeel, P. (2016). *End-to-End Training of Deep Visuomotor Policies*. JMLR, 17(1), 1334-1373.
- Chi, C. et al. (2023). *Diffusion Policy: Visuomotor Policy Learning via Action Diffusion*. arXiv:2303.04137.
- Zhao, T. Z., Kumar, V., Levine, S., Finn, C. (2023). *Learning Fine-Grained Bimanual Manipulation with Low-Cost Hardware* (introduces ACT). arXiv:2304.13705.