# Long-horizon manipulation

Hard problem for RL

- **High precision** tasks

- **Complex reward** design

- **Large exploration** space

RL becomes too **inefficient**

# Multi-stage feedback

Long-horizon tasks have a **multi-stage structure**

**Key idea:**
**Stage feedback** + **demos** can guide learning

**DEMO3:**
**Sample efficient** RL for **multi-stage manipulation**



Stage 1: Grasp

Stage 1: Grasp

Stage 2: Align with hole
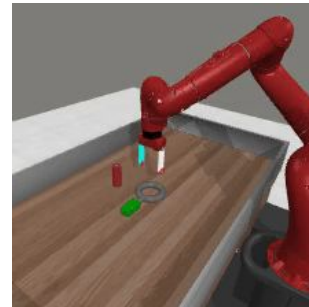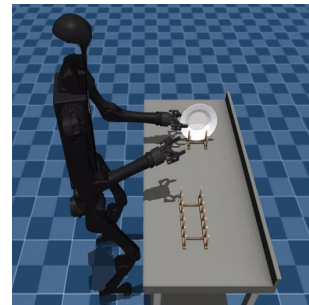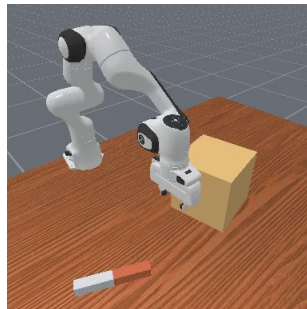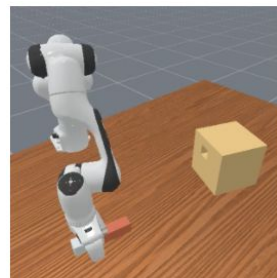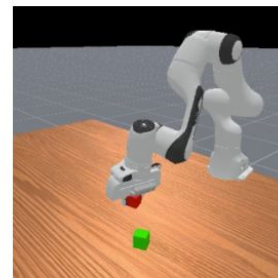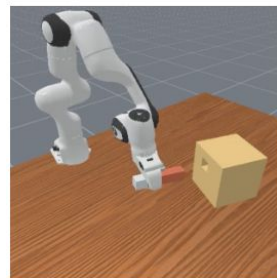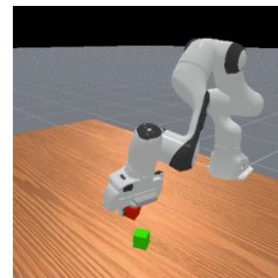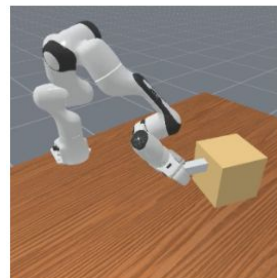
Stage 2: Hover over cube

Stage 3: Insert

Stage 3: Stack

**DEMO³**: Multi-Stage Manipulation with Demonstration-Augmented **Reward**, **Policy** and **World-Model** Learning

**How do we achieve this?**

**Demonstration to simultaneously learn:**

**1. Online reward function**

**2. Policy pre-training**

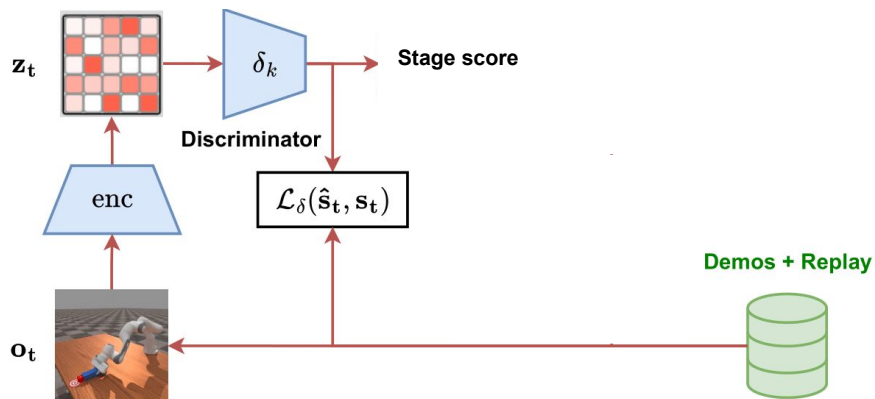**3. World Model for planning**

How do we achieve this?

Demonstration to simu...

From images observations!

1. Online reward function

2. Policy pre-training

3. World Model for planning

**DEMO³**: Multi-Stage Manipulation with Demonstration-Augmented **Reward**, **Policy** and **World-Model** Learning
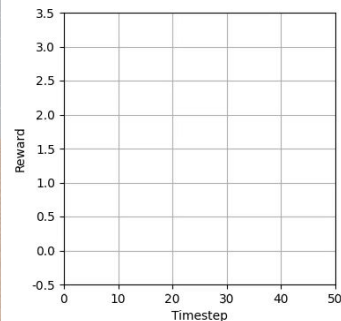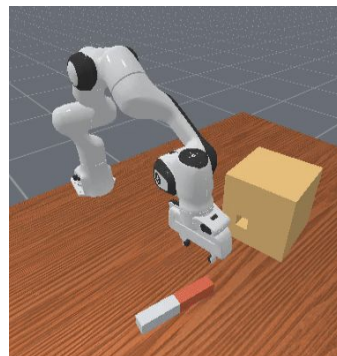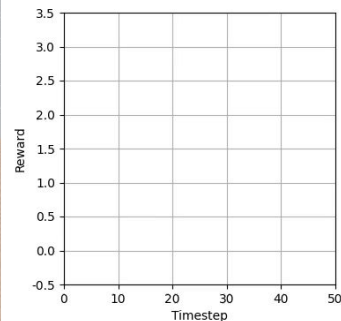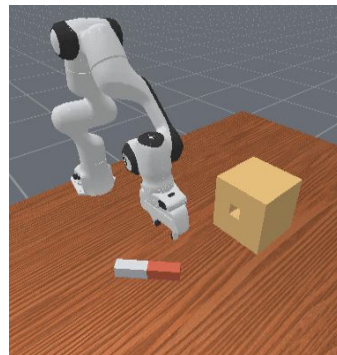
# Online reward learning

Use demos for **online reward learning**:

- Dataset: Demos + Replays

- Discriminator separates success / fail frames
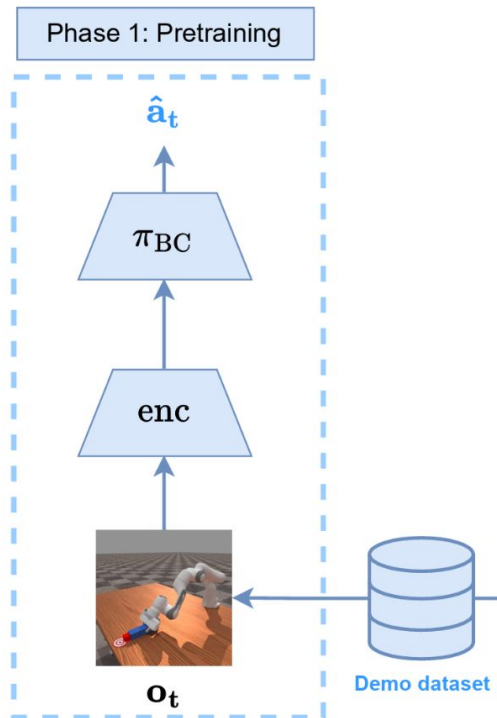
- Policy trained on learned dense reward

$z_t$

$\delta_k$    Stage score

**Discriminator**

enc

$\mathcal{L}_\delta(\hat{s}_t, s_t)$

$o_t$

**Demos + Replay**

**Successful rollout**

**Failed rollout**

# Policy pre-training

Phase 1: Pretraining

$\hat{a}_t$
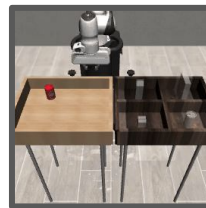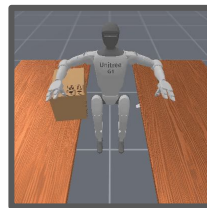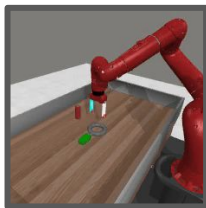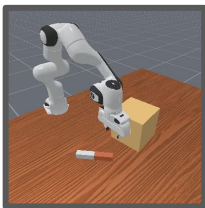
$\pi_{BC}$

enc

$o_t$

Demo dataset

Behavioral cloning warm-starts **policy** and **encoder**

We sample from BC policy in **early stages**

**Progressively transition** to RL policy

# World Model Learning

**DEMO³**: Multi-Stage Manipulation with Demonstration-Augmented **Reward**, **Policy** and **World-Model** Learning
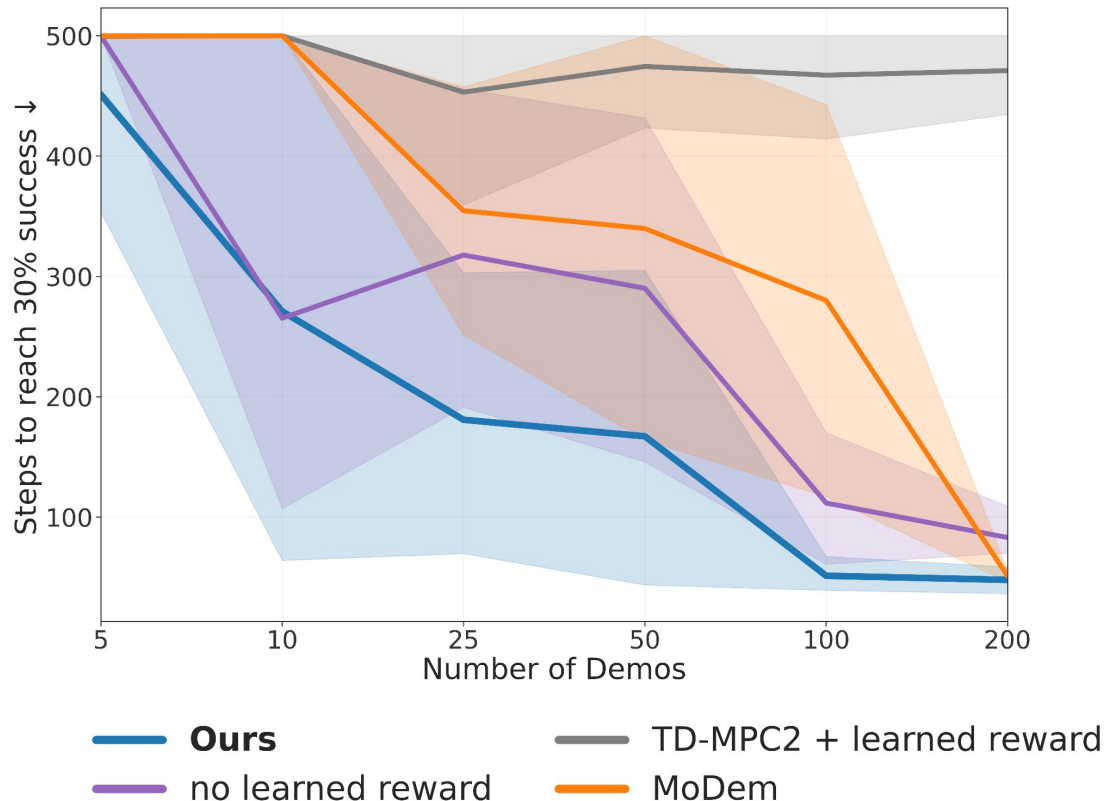
# Results: Sample Efficiency



We solve hard tasks in < 100K Steps from **only images and sparse rewards**

# Results: Demonstration Efficiency

DEMO3 excels at **demonstration efficiency**

Hard tasks solved with **only 5 demos**

# Thank you!