



**ICLR**  
International Conference On  
Learning Representations

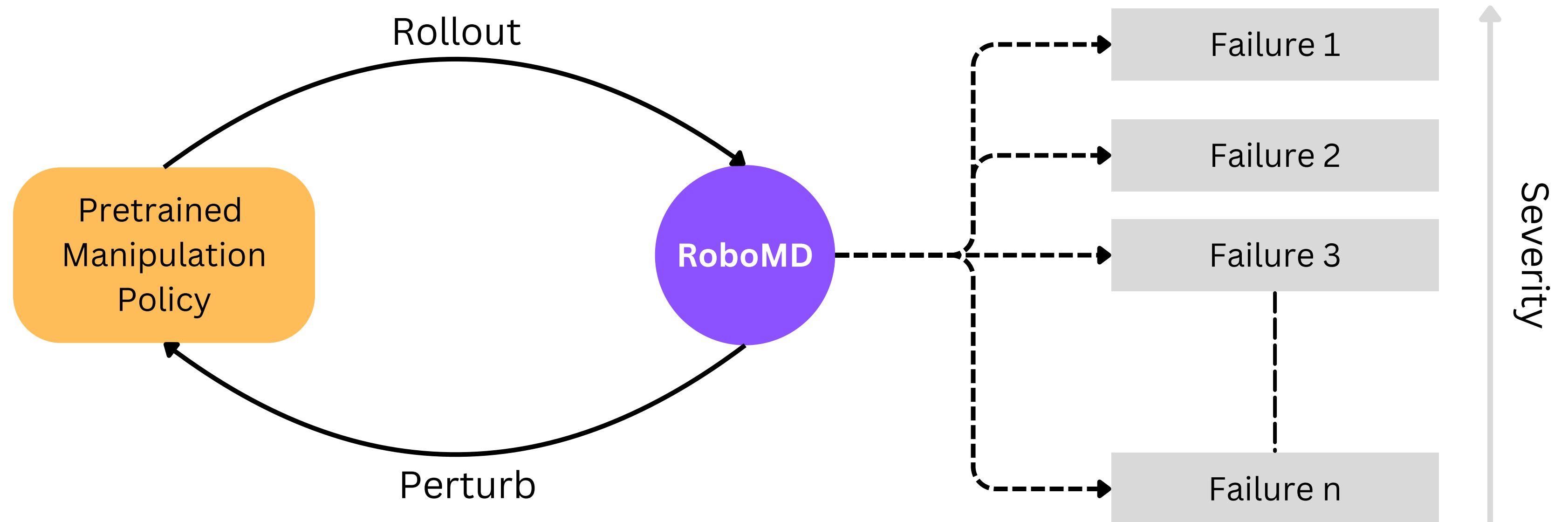
# RoboMD: Uncovering Robot Vulnerabilities through Semantic Potential Fields

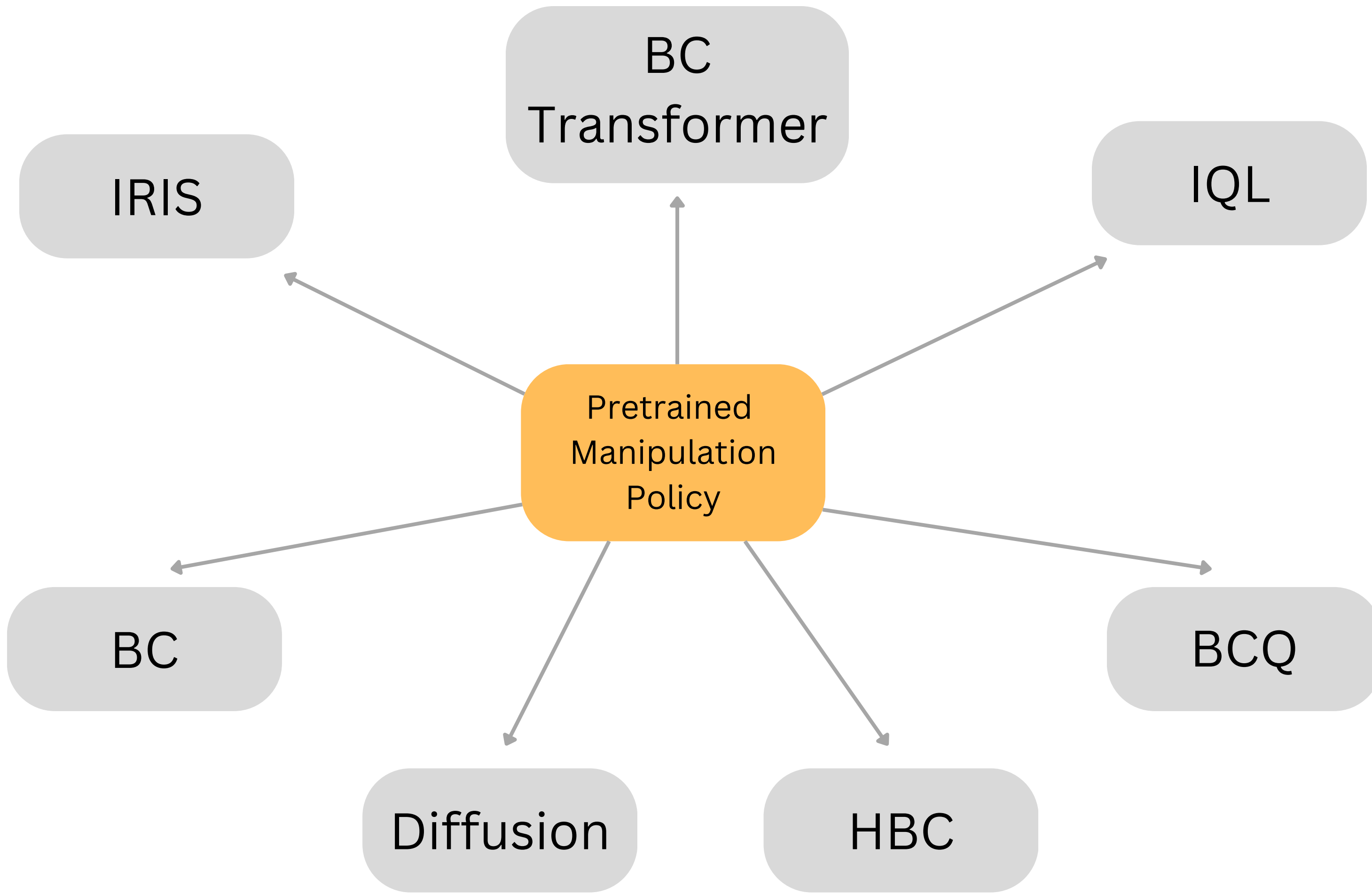
---

Som Sagar, Jiafei Duan, Sreevishakh Vasudevan, Yifan Zhou,  
Heni Ben Amor, Dieter Fox, and Ransalu Senanayake



Our framework, **RoboMD**, provides a systematic approach for **diagnosing and quantifying failure modes** in **robot manipulation policies**, enabling robust performance improvement across **diverse and unseen environments**

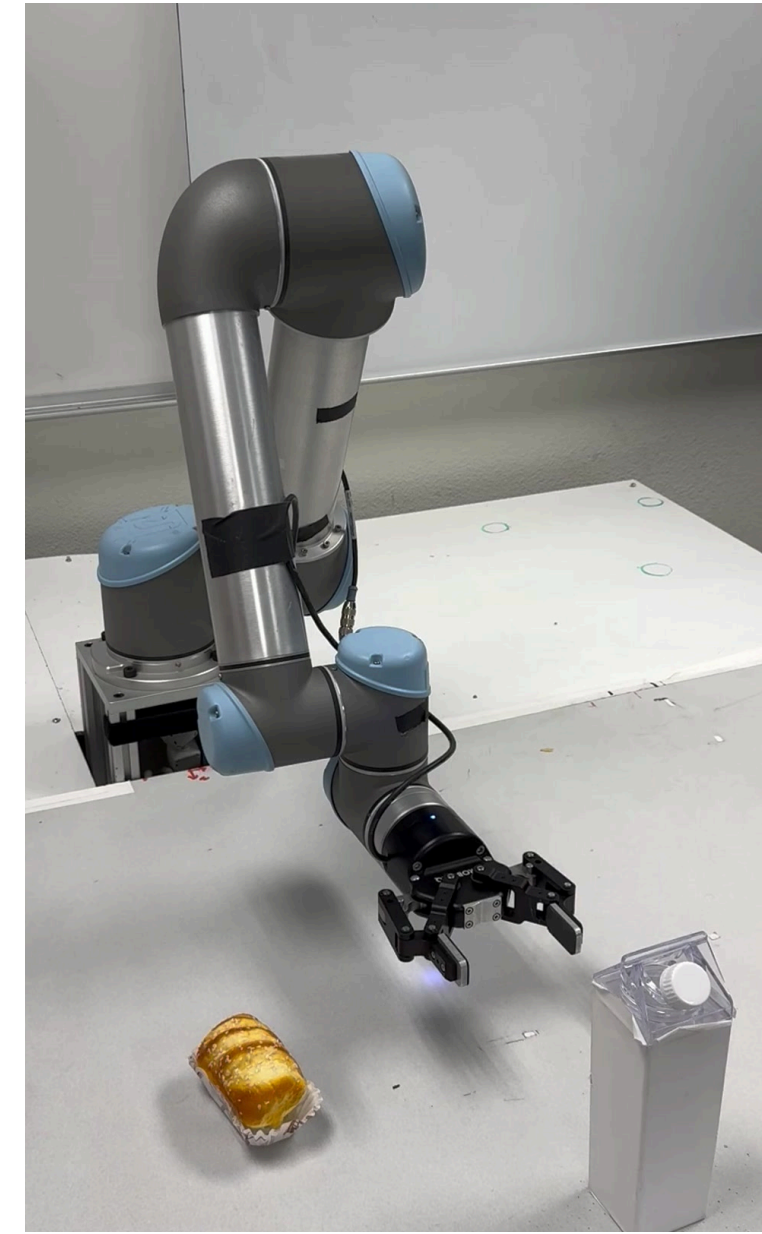
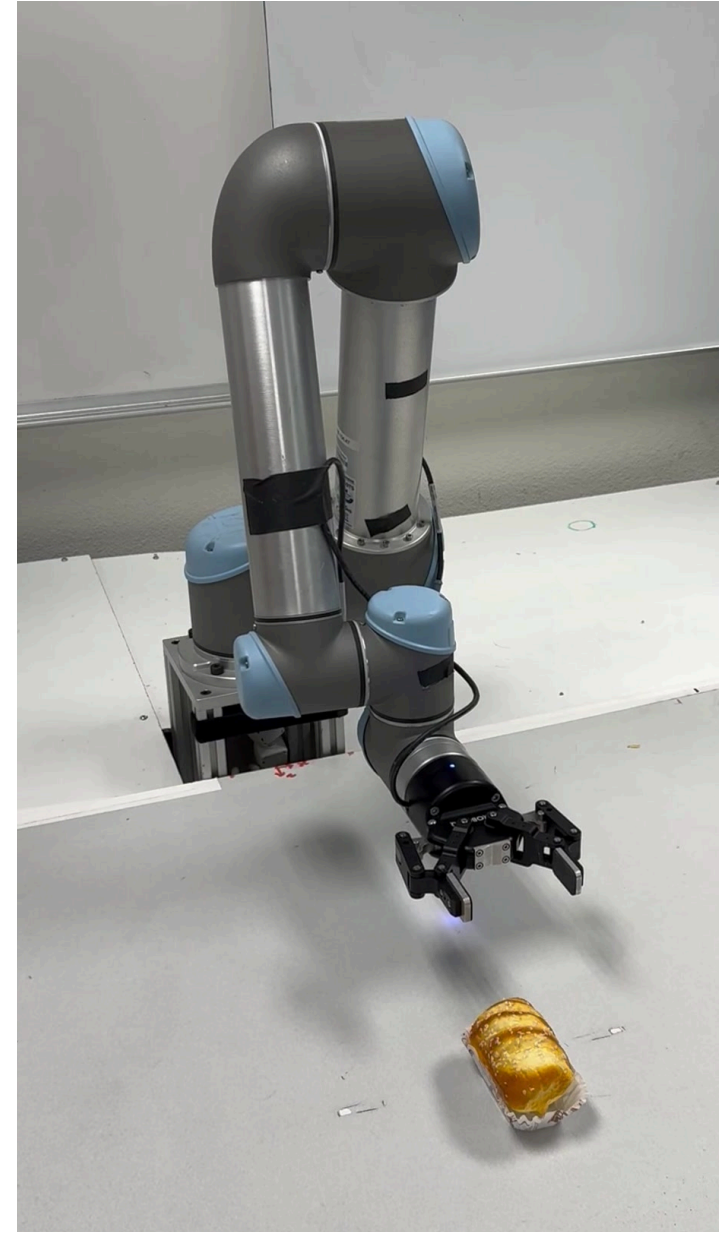
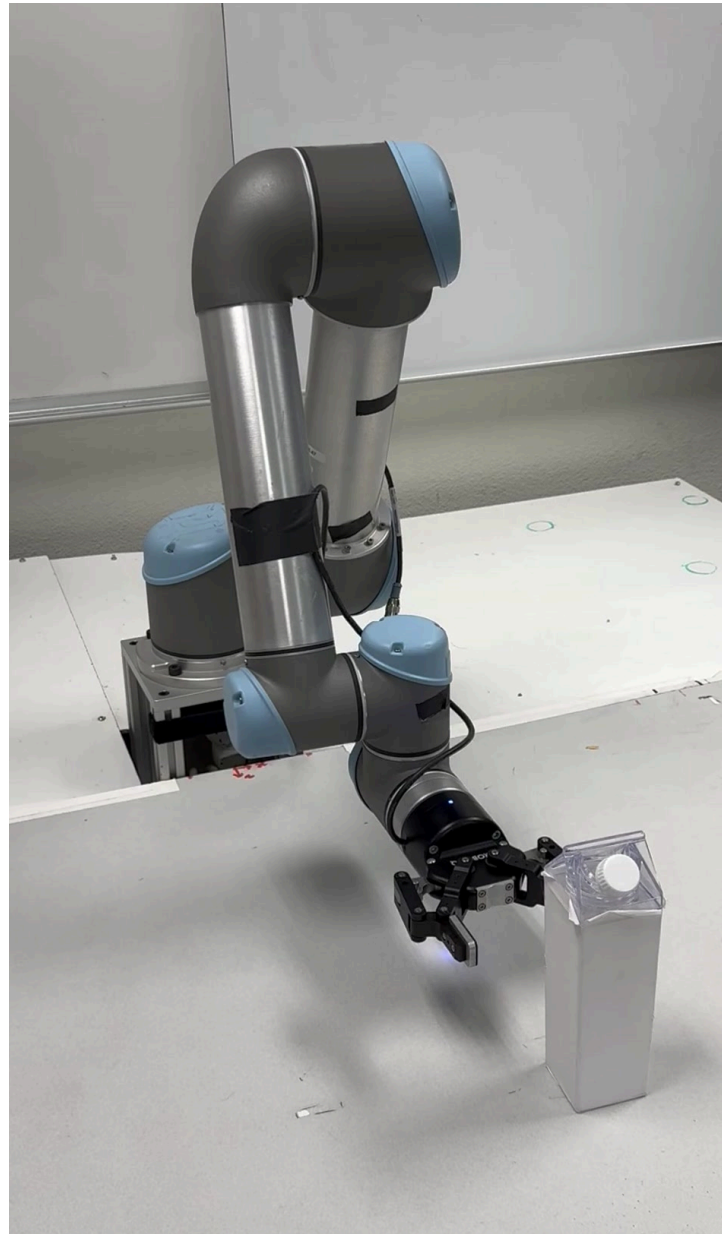




Pretrained  
Manipulation  
Policy

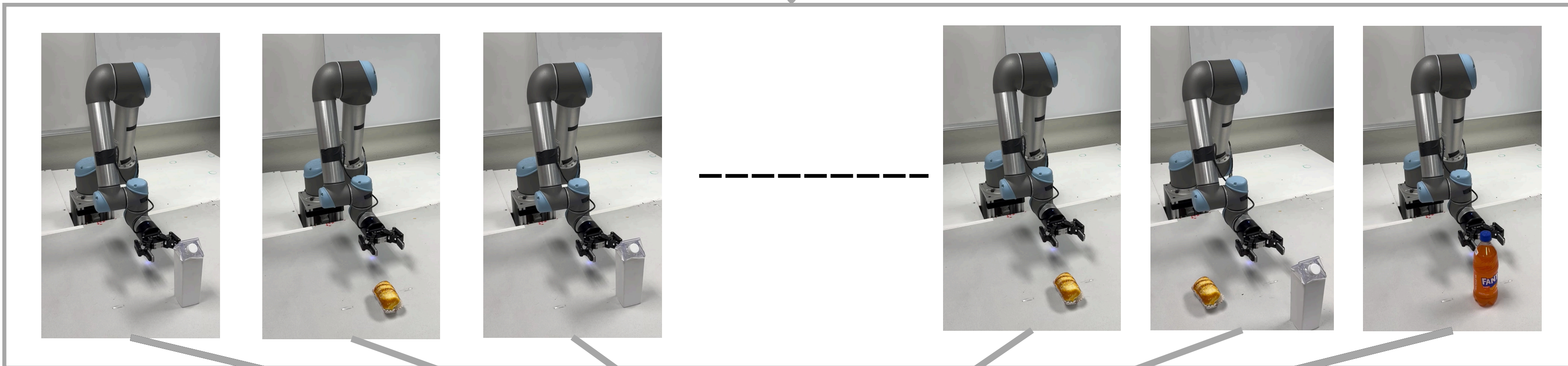


Rollout



Pretrained Manipulation Policy

Choose rollout of Milk Carton



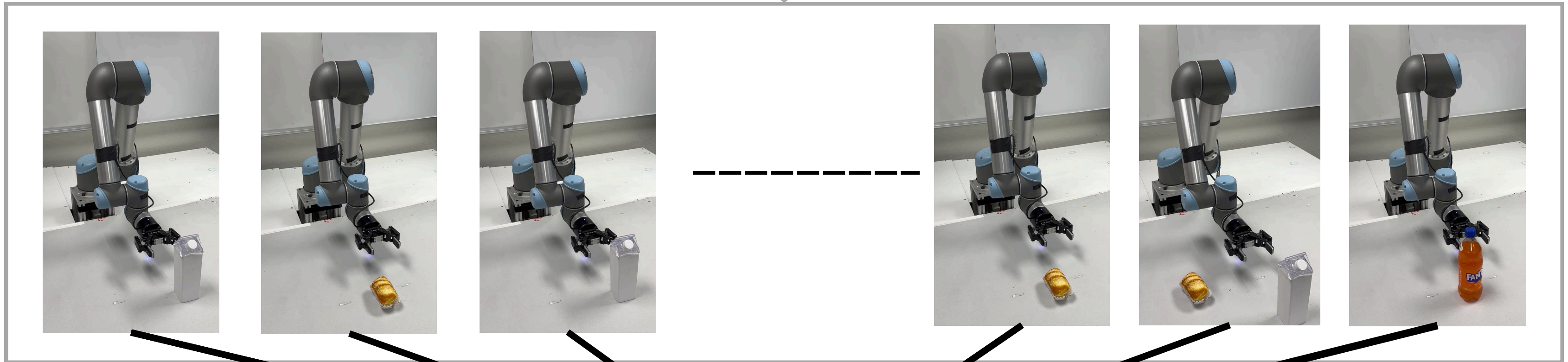
Action "Milk Carton"

Training RoboMD

Reward based on failure during rollout

Pretrained  
Manipulation  
Policy

Through **RL RoboMD** observes multiple observations and learns a distribution of **failure over multiple actions (Bread, Fanta, Milk Carton)**



Training  
RoboMD

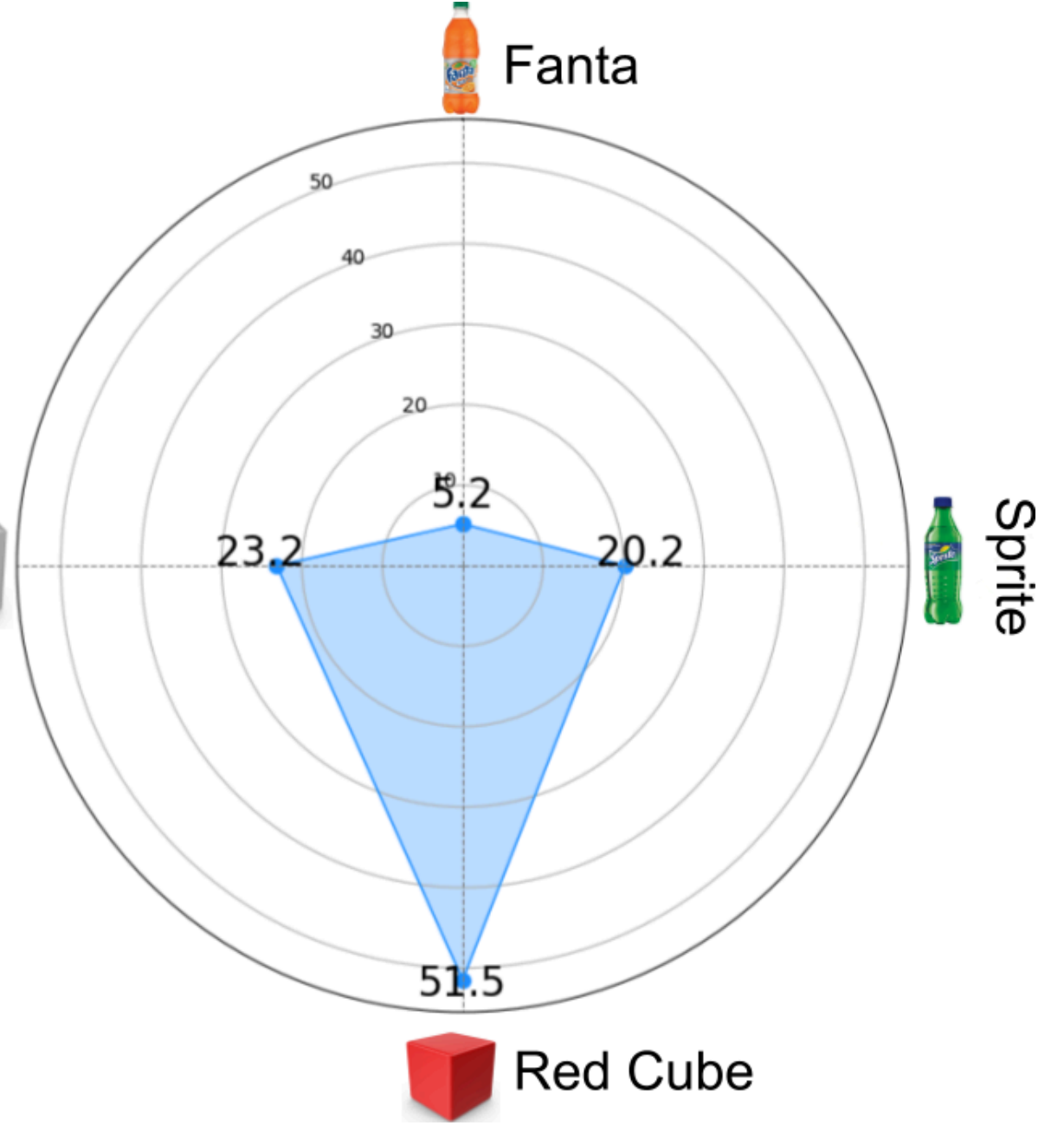
Pretrained  
Manipulation  
Policy



Trained  
RoboMD



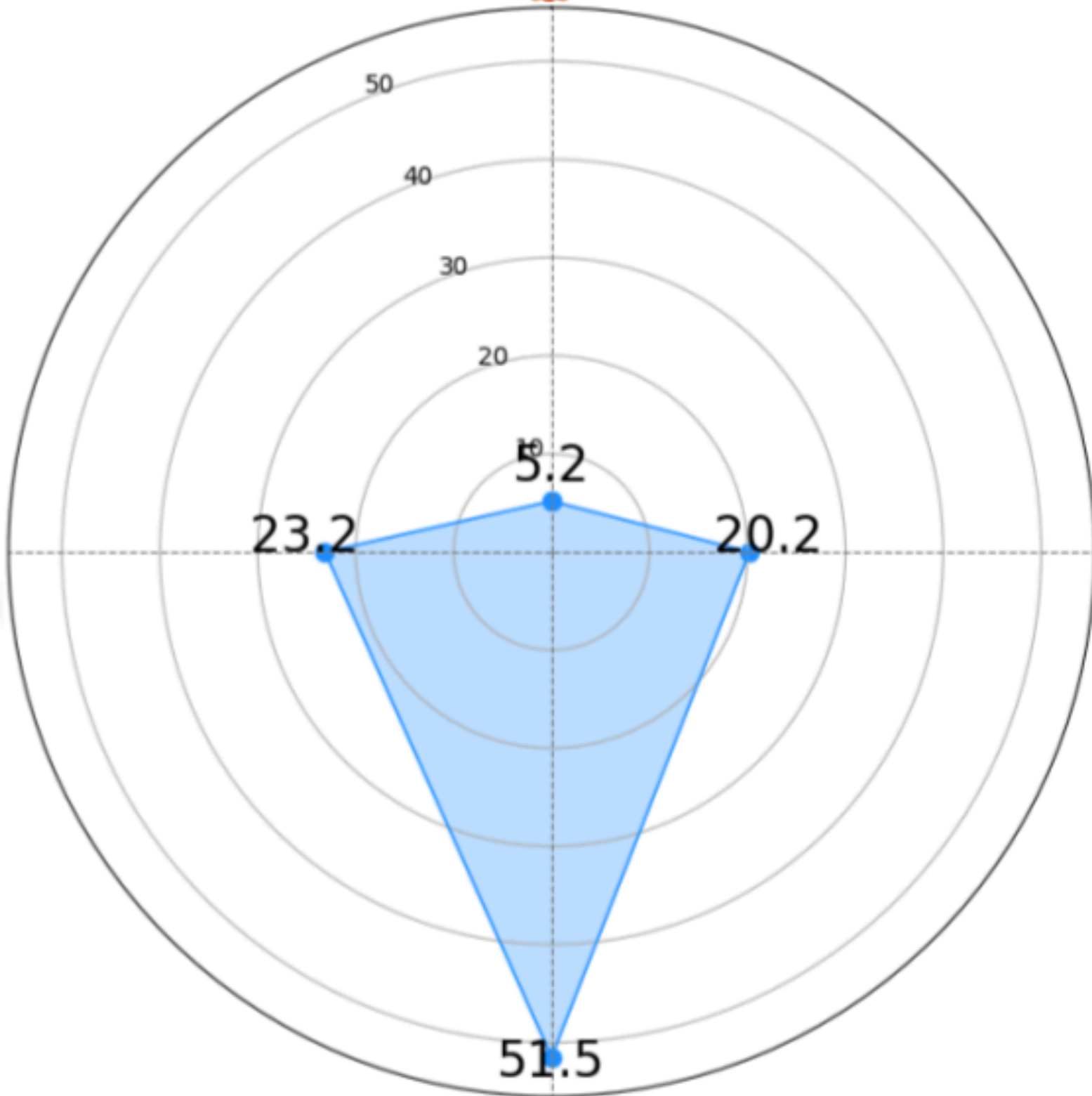
Milk Carton



Trained RoboMD



Milk Carton



Fanta



Sprite

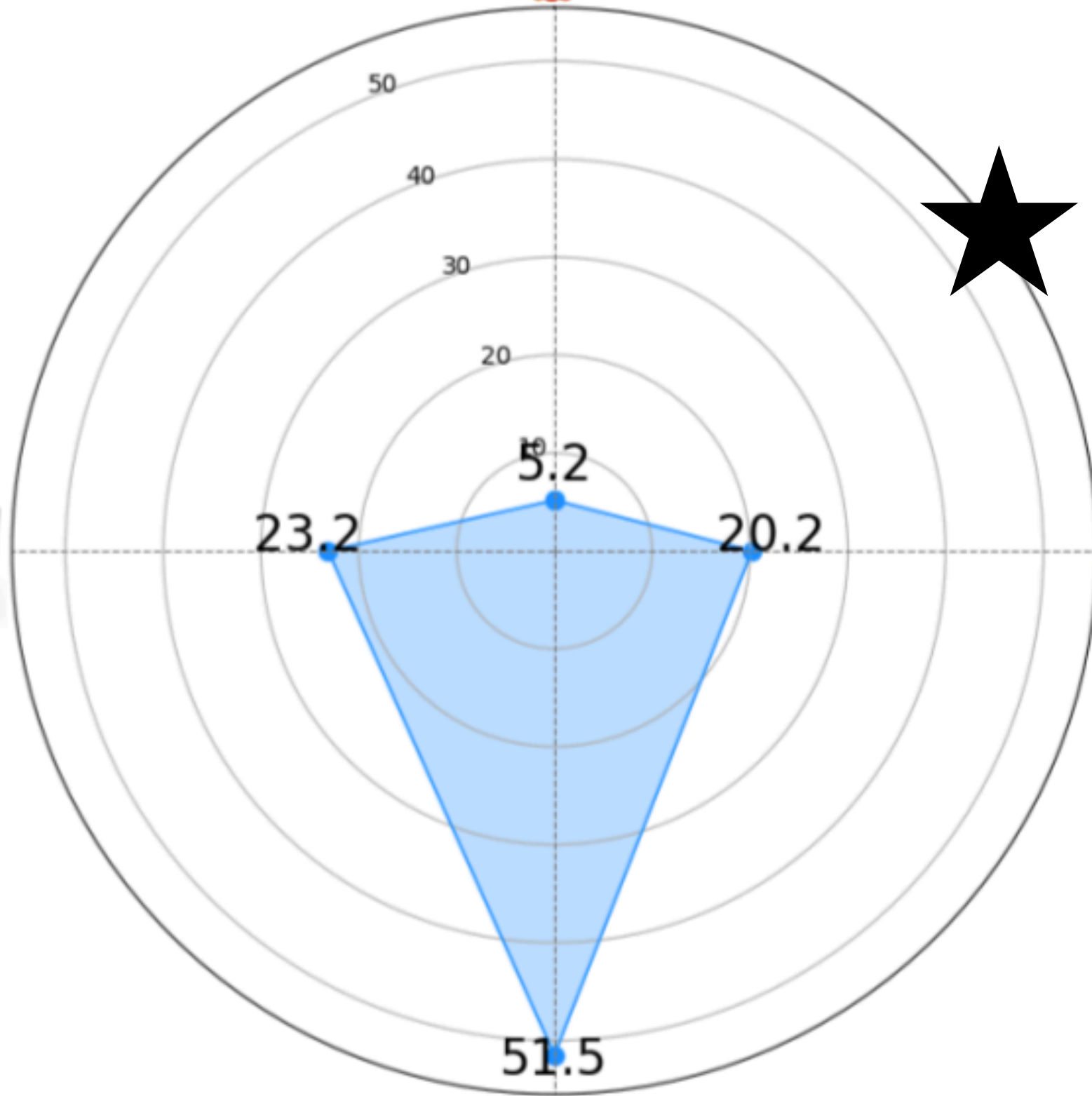


Red Cube

Trained RoboMD



Milk Carton



Fanta



Sprite

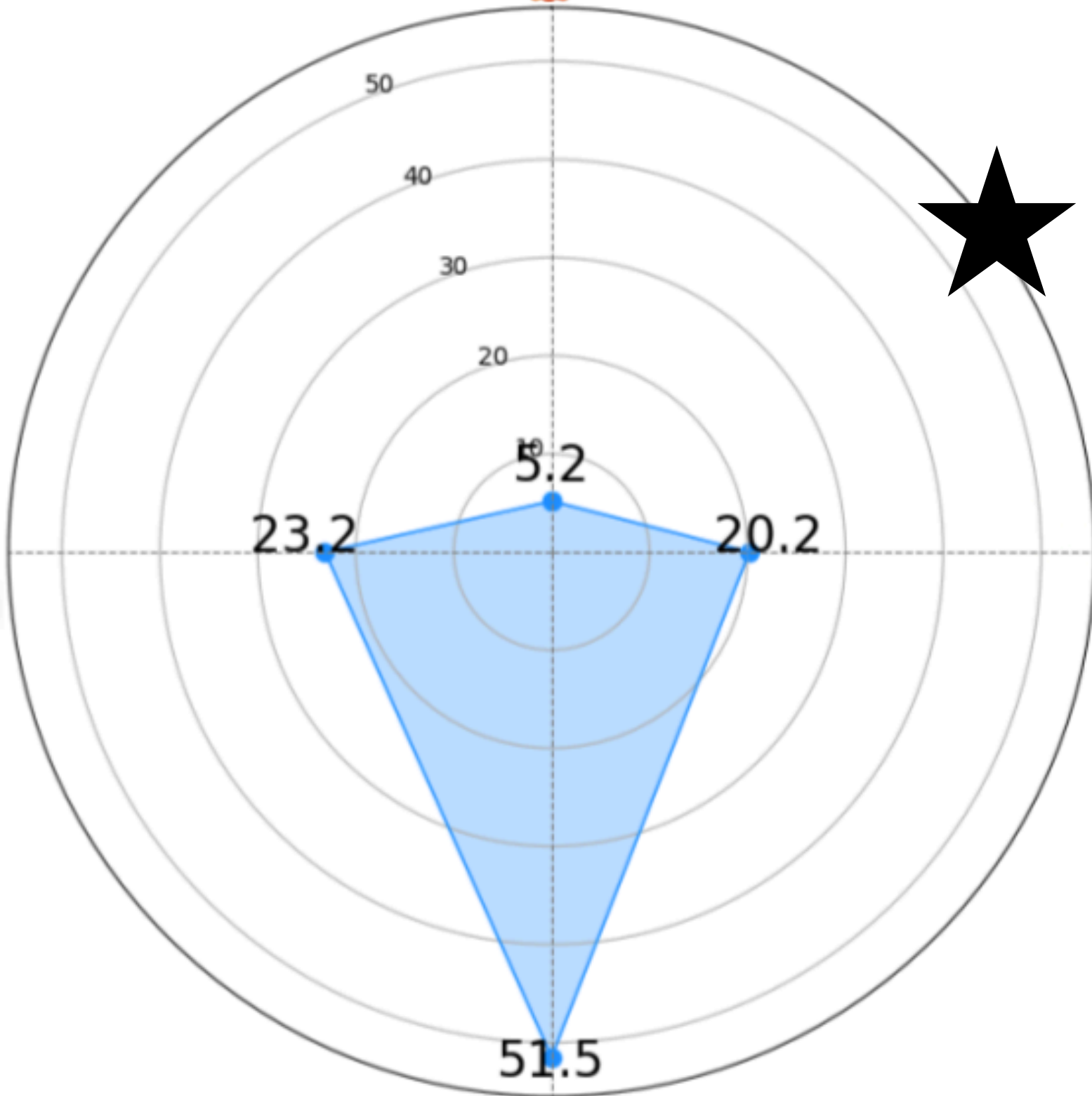


Red Cube

Trained  
RoboMD



Milk Carton



??



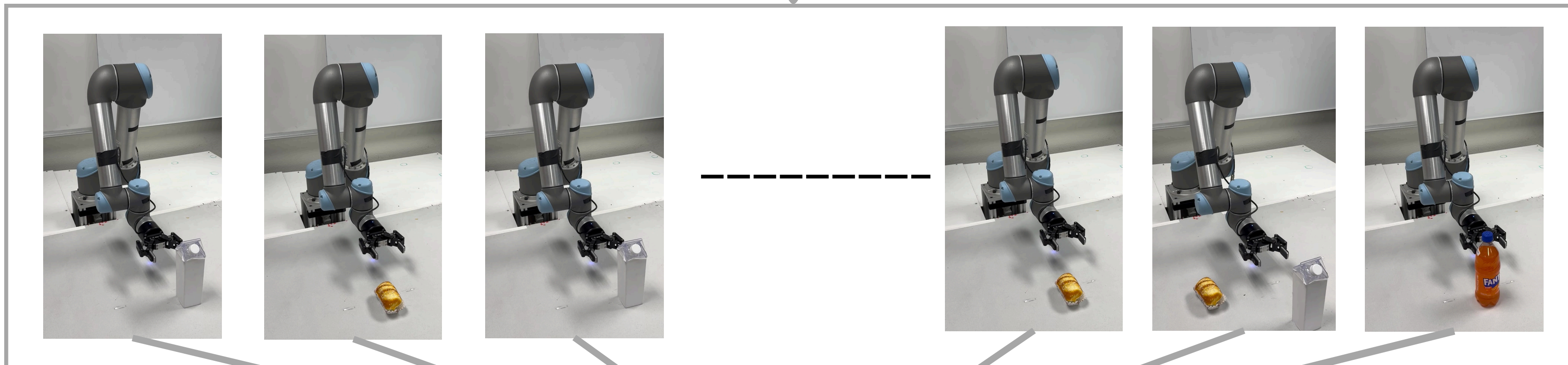
Sprite



Red Cube

Pretrained  
Manipulation  
Policy

Change RL action space to  
handle continuous action space

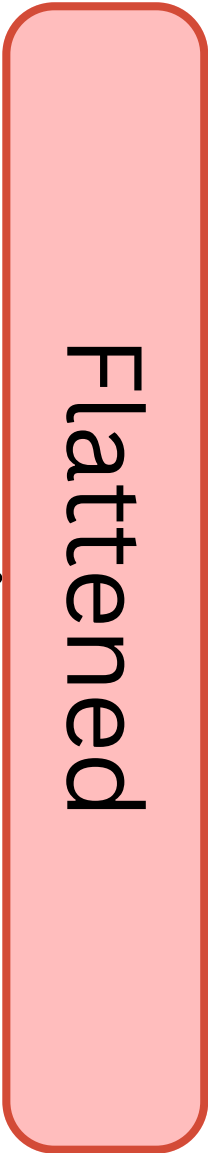
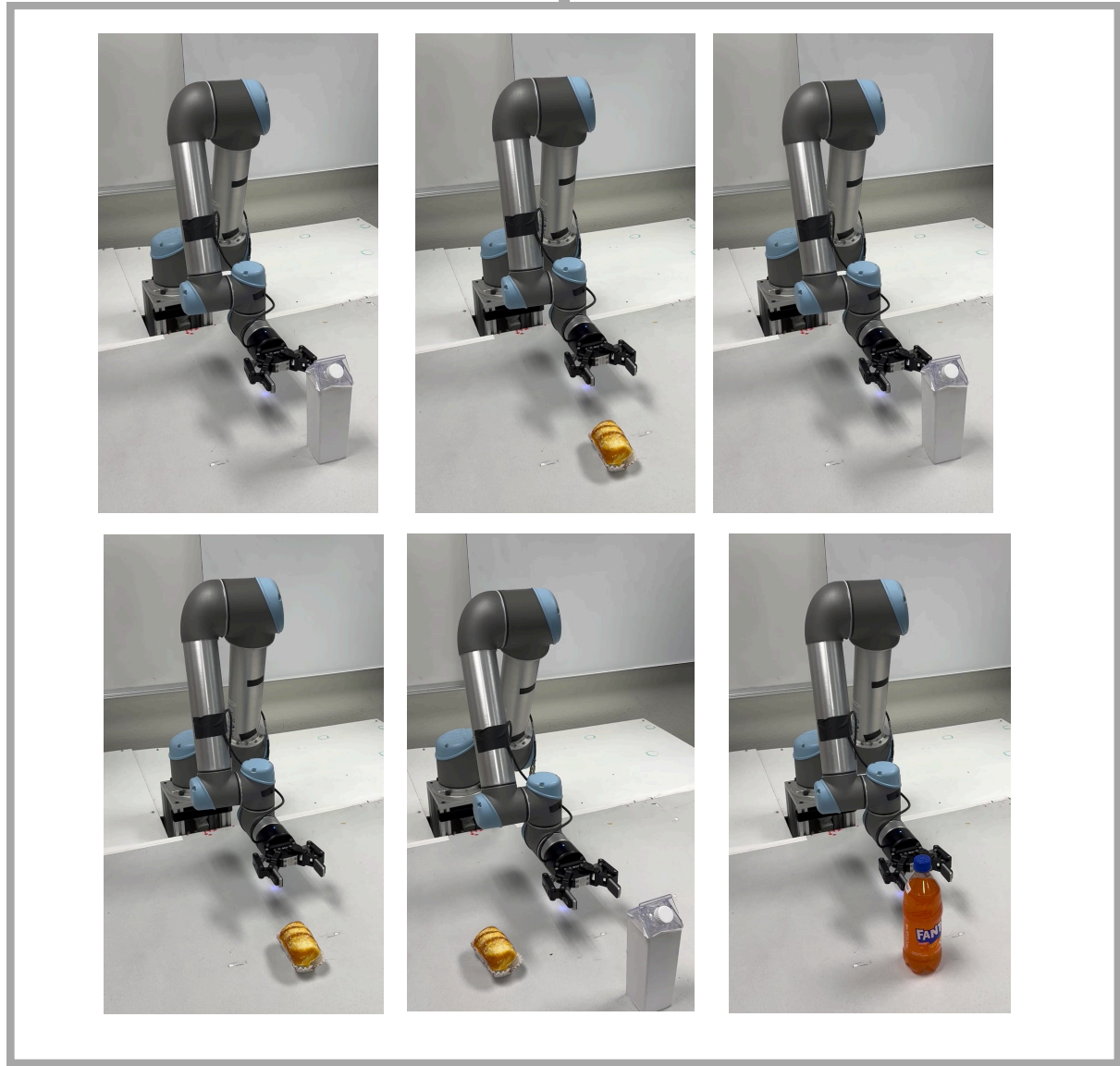


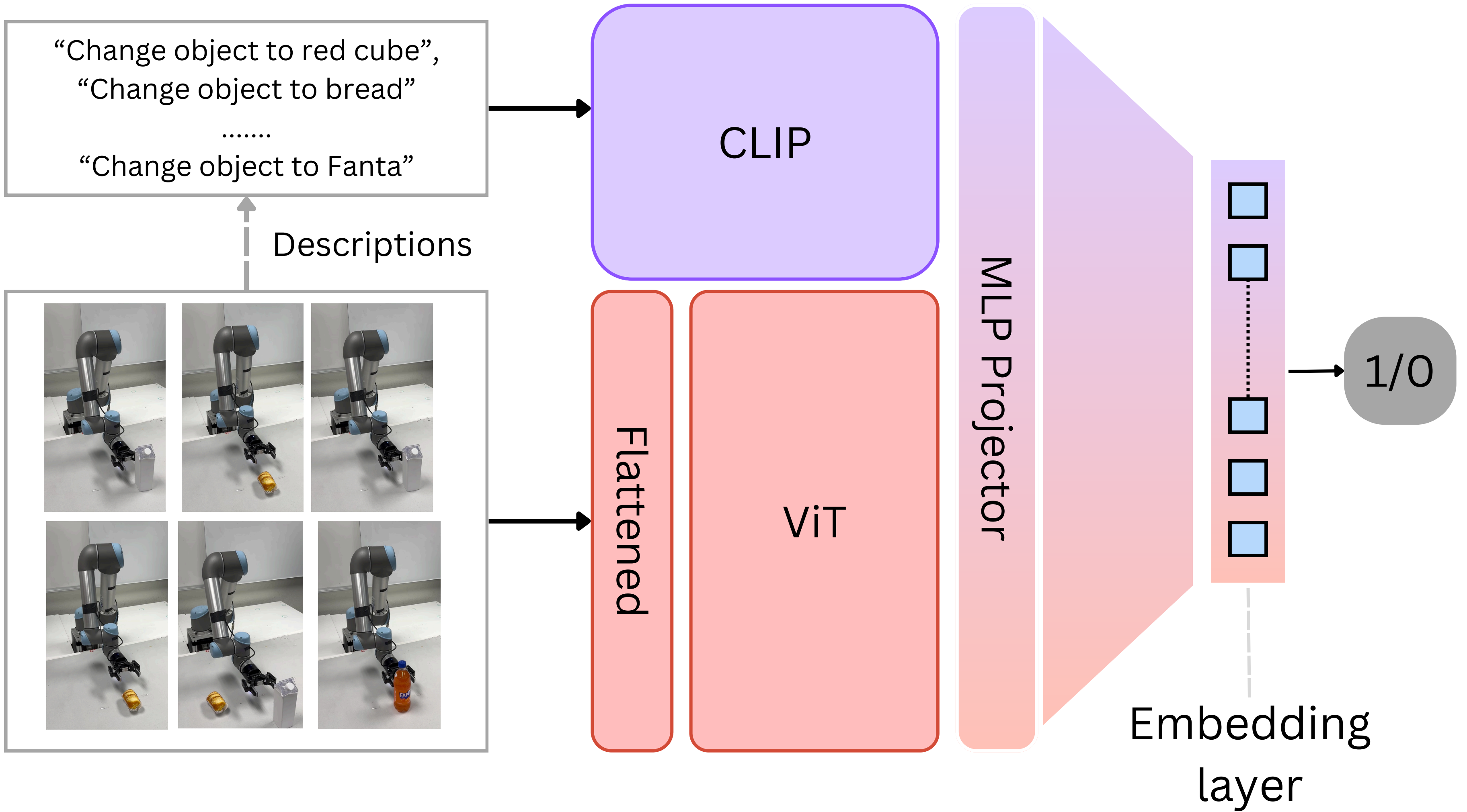
Training  
RoboMD

“Change object to red cube”,  
“Change object to bread”  
.....  
“Change object to Fanta”

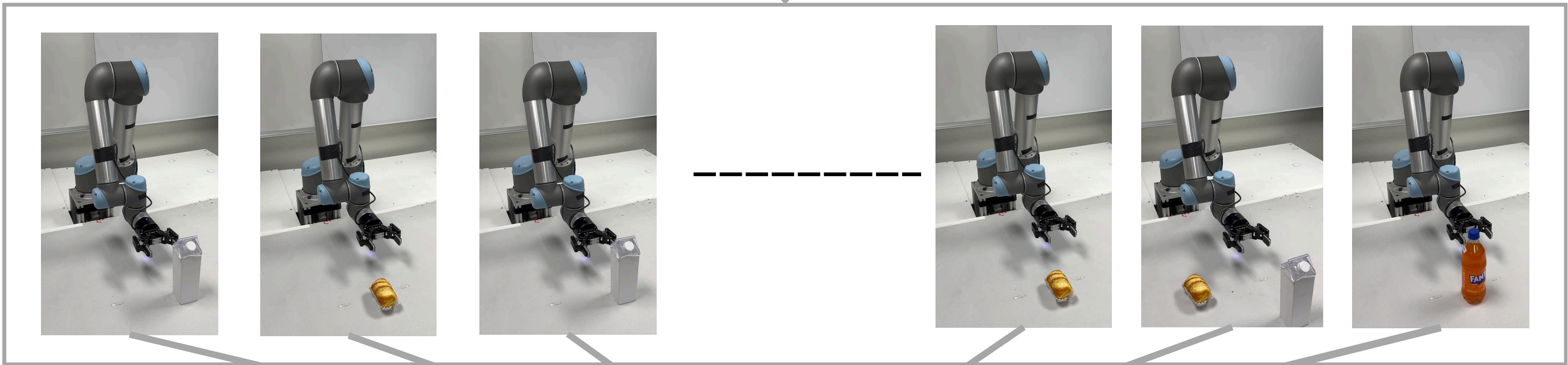


Descriptions





Pretrained  
Manipulation  
Policy

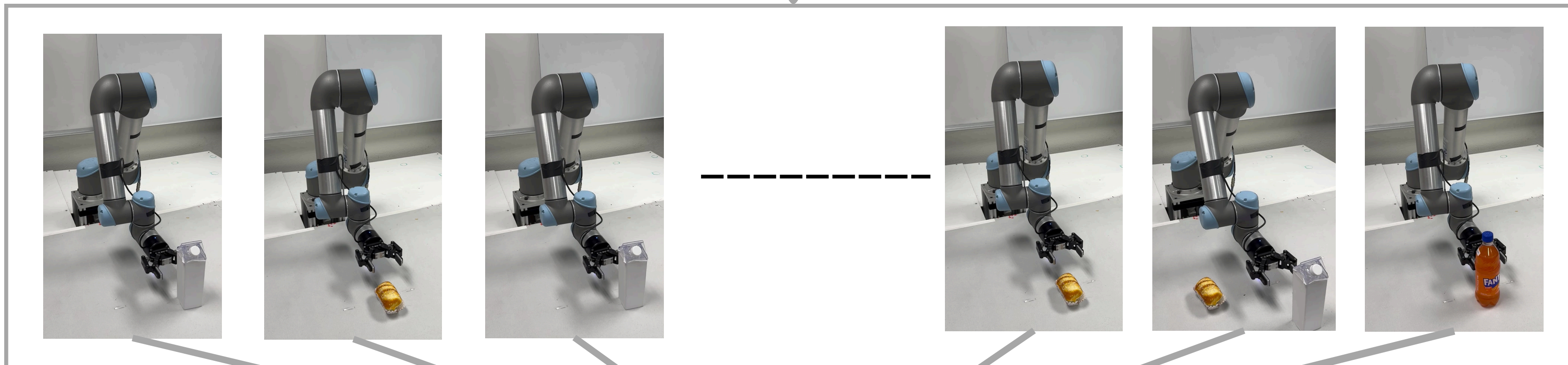


Training  
RoboMD

Search for closest known embedding

Pretrained Manipulation Policy

Choose rollout of Fanta



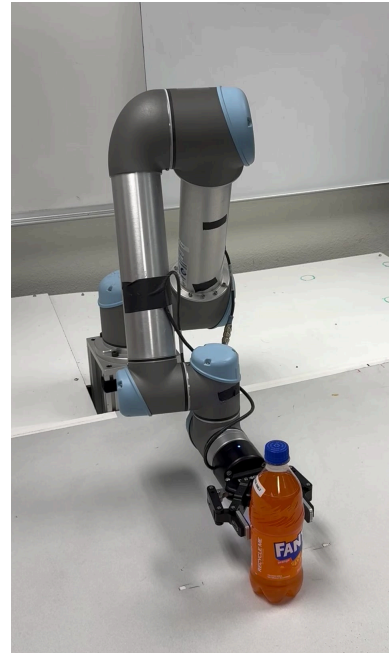
Random action embedding

Training RoboMD

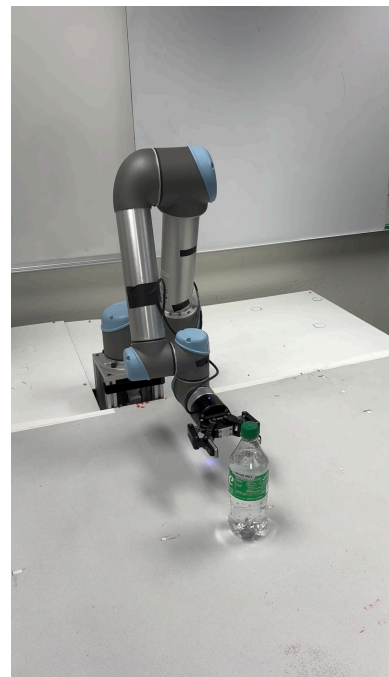
$$\text{Reward} \propto \frac{1}{d(e_{\text{known}}, e_{\text{unknown}})}$$

Rank seen and unseen scenarios  
based on their likelihood to fail!

Rank seen and unseen scenarios based on their likelihood to fail!



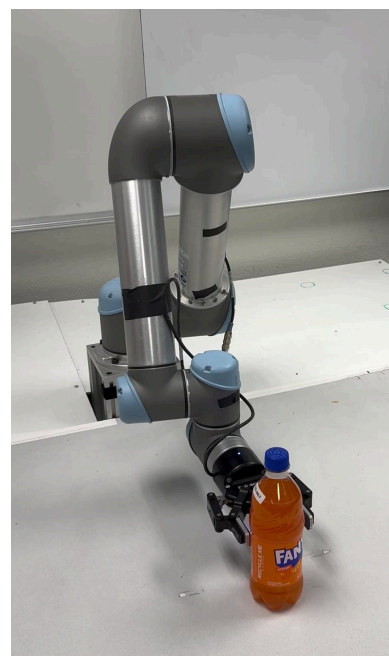
Seen



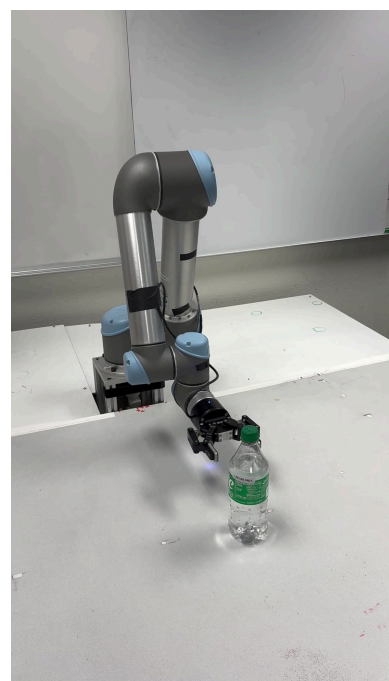
Unseen

Trained  
RoboMD

Rank seen and unseen scenarios based on their likelihood to fail!



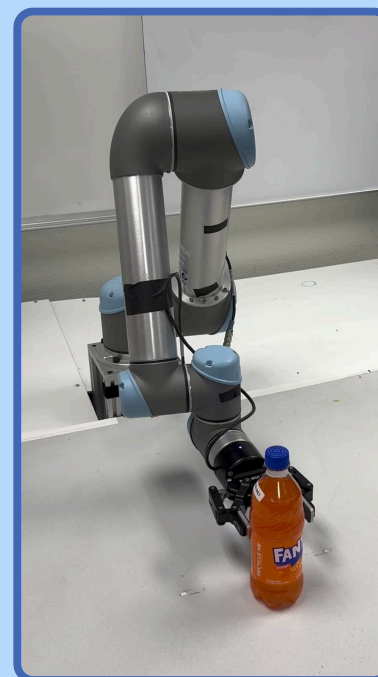
Seen



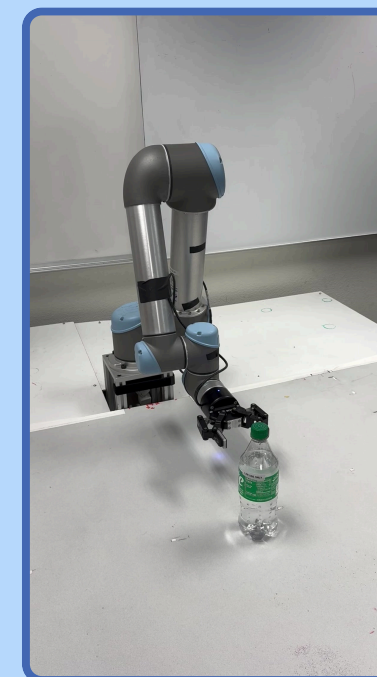
Unseen

Trained  
RoboMD

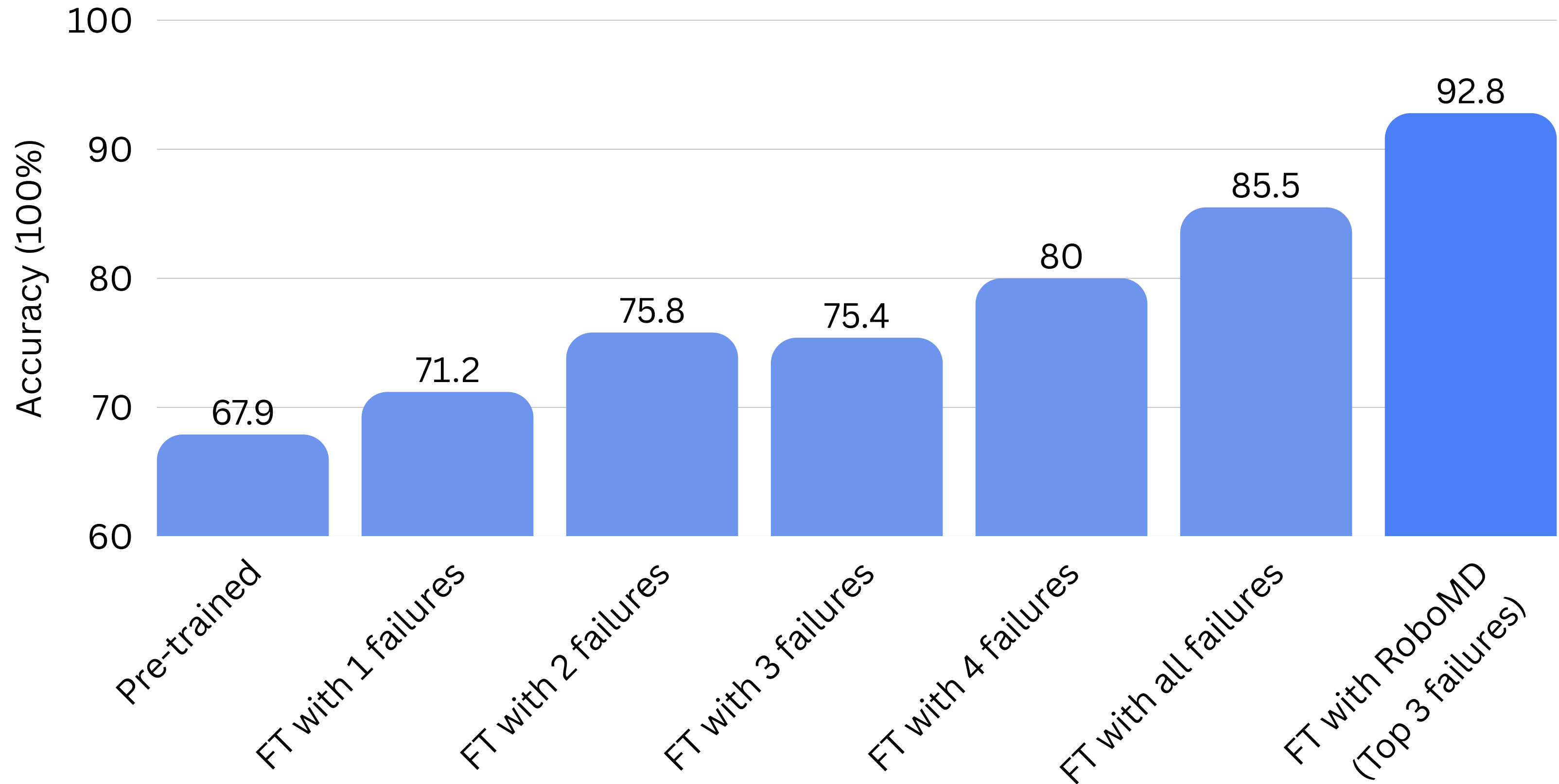
The manipulation policy is more likely to fail with Fanta bottle than water bottle.



>



# Improve the Manipulation Policy with RoboMD





Paper



Github