

### Step 1: Pose prediction

*Initial state:*



*Language prompts:*

*Predict the contact point and orientation for pulling the {object}*

Input



MCC-MLLM

*Pose prediction:*

Contact point:  $[x, y]$ ,  
gripper's upward direction:  $[x_w, y_w, z_u]$ ,  
gripper's forward direction:  $[x_f, y_f, z_f]$ .

### Step 2: Failure detection and correction

*End state:*



*Language prompts:*

*The robot's end-effector state is .....  
Detect the failure causes of pulling the {object}.*

*Correction experts:*



Position expert



Reasoning expert



Rotation expert



Input



Manipulation and Correction Collaborative (MCC)-MLLM

Error type

Prompt feedbacks

*Failure causes:*

*The failure cause is  
{position // rotation //  
position and rotation}.*

*Corrected pose prediction:*

Contact point:  $(x^c, y^c)$ ,  
gripper's upward .....

### Step 3: Continuous policy learning

*Corrected state:*



*Corrected pose label:*

Contact point:  $[x^c, y^c]$ ,  
gripper's upward .....

Train samples



MCC-MLLM

Step 3

Step 2

Step 1

Close-loop  
Correction



: Fine-tune



: Frozen