



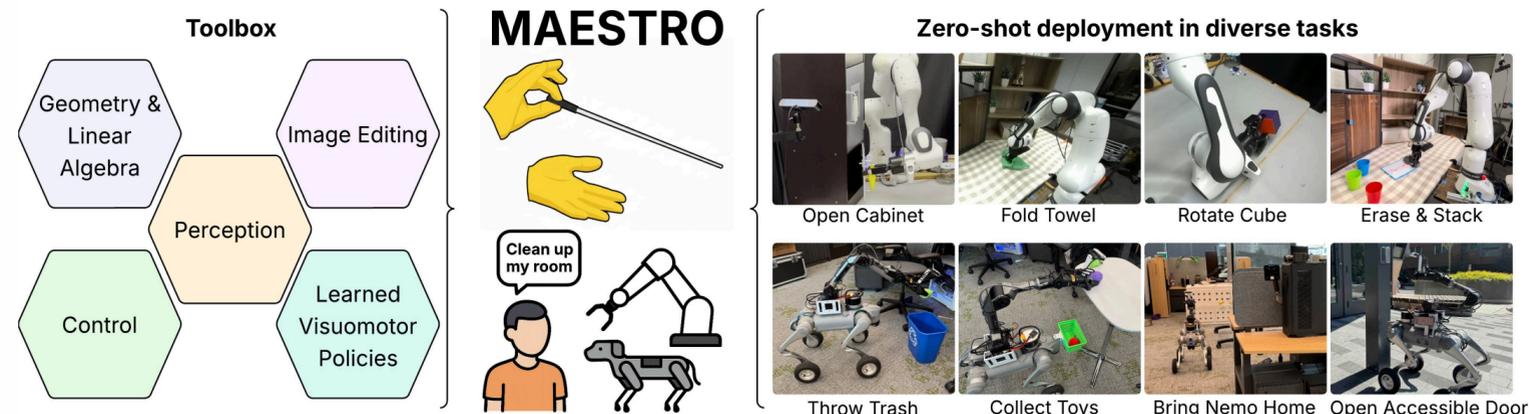
# Maestro

## Orchestrating Robotics Modules with Vision-Language Models for Zero-Shot Generalist Robots



Junyao Shi\*, Rujia Yang\*, Kaitian Chao\*, Bingqing Wan, Yifei Shao, Jiahui Lei, Jianing Qian, Long Le, Pratik Chaudhari, Kostas Daniilidis, Chuan Wen, Dinesh Jayaraman

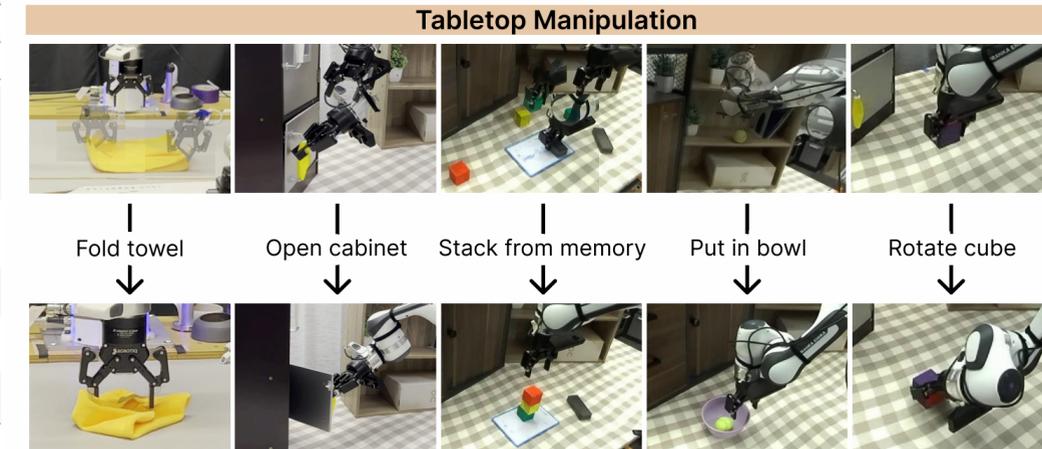
### Overview



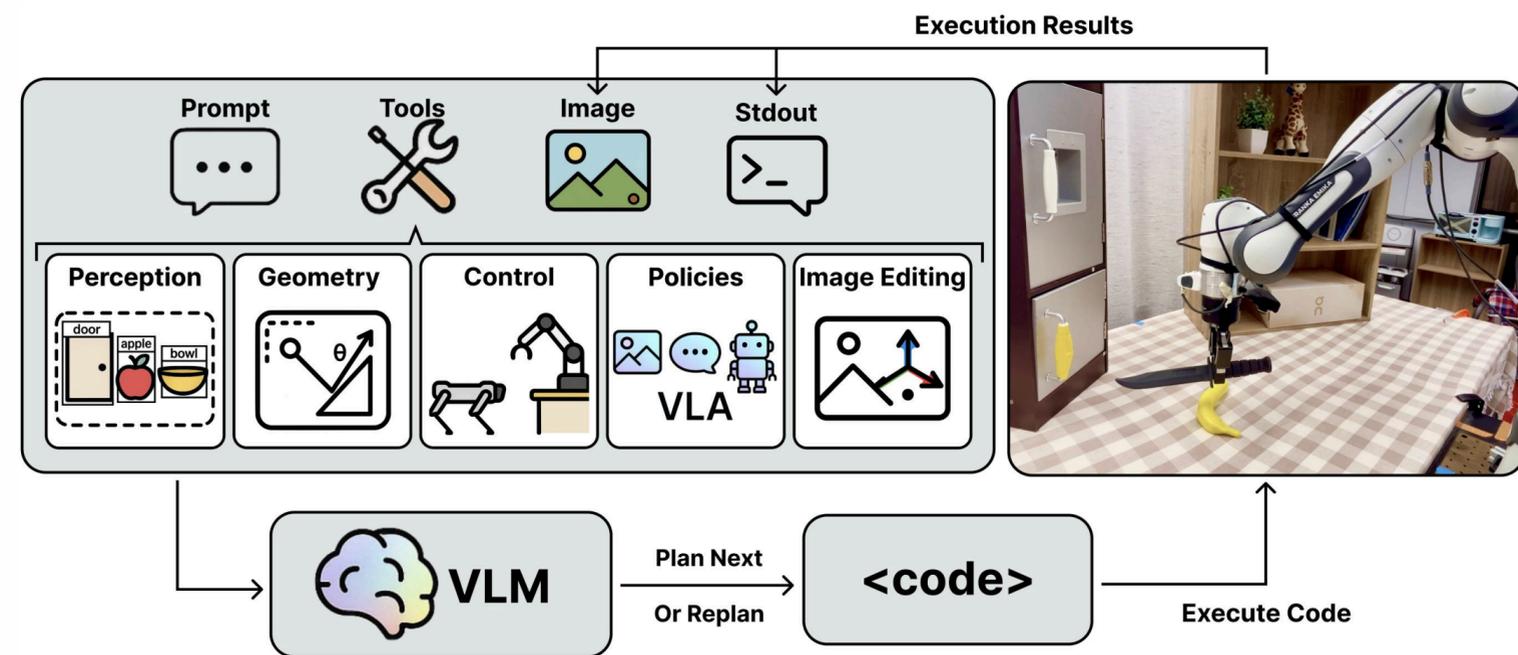
### System Tools

Tool Category	Prior Work Examples	MAESTRO Modules
<b>Tabletop Manipulation</b>		
Perception	Raw sensory inputs (RGB + proprioception); Segmentation/Bounding Box centers [2, 8, 13]	Raw sensory inputs (RGB + proprioception); Segmentation centers; <b>Active perception (zoom/look around with wrist camera)</b> ; FoundationStereo [25] depth; Gemini pointing; <b>VLM-selected task-relevant keypoints (ReKep-inspired [26])</b>
Control	Cartesian control, gripper control [2, 8]; Movement primitives [13]	Cartesian control, gripper control; <b>cuRobo collision-free motion planning</b>
Learned Visuomotor Policies	m2t2 grasp model [27]	<b>GraspGen grasp model [28]; <math>\pi_{0.5}</math> VLA with high-frequency closed-loop monitoring (Qwen-2.5-VL)</b>
Geometry & Linear Algebra (new)	<i>None</i>	<b>Distance measurement, vector construction, vector rotation, relative rotation between vectors</b>
Image Editing (new)	<i>None</i>	<b>Draw points, overlay 6D poses to improve visual grounding</b>
<b>Extra Modules for Mobile Manipulation</b>		
Perception	Build global/local map [6]	Mobile base state estimation; <b>Active perception tools (look left/right/ground, view carry-on basket, log object location)</b>
Locomotion	Navigation [6]	Navigation; <b>Fine-grained "nudge" tool for local adjustment</b>

### Experiments



### Method



### Results

<b>Tabletop Manipulation</b>					
Challenge	Task Description	Gemini Robotics Agent	$\pi_0$	$\pi_{0.5}$	MAESTRO
Pick-Place	Put item in bowl	73.3 ± 46.2	74.0 ± 37.1	70.0 ± 41.1	<b>98.0 ± 4.5</b>
Deformable object	Fold the four corners of the towel into the center	40.0 ± 17.3	47.0 ± 25.1	70.0 ± 15.4	<b>71.3 ± 21.4</b>
Articulated object	Open cabinet	3.3 ± 5.8	8.3 ± 2.9	0.0 ± 0.0	<b>68.0 ± 31.3</b>
Spatial reasoning	Rotate cube purple side up	23.6 ± 3.5	29.0 ± 1.7	10.0 ± 0.0	<b>60.0 ± 38.1</b>
Memory & long-horizon & semantic	Erase instructions on whiteboard, then follow instruction to stack cups	26.7 ± 24.7	12.0 ± 12.0	22.0 ± 22.8	<b>63.0 ± 16.8</b>

<b>Mobile Manipulation</b>		
Task Category	Task Description	MAESTRO
Long-horizon manipulation	Collect all toys on table	85.0 ± 22.4
Long-horizon loco-manipulation	Throw green ball into garbage can	76.7 ± 14.9
Active exploration	Search item and put on table	96.0 ± 8.9
Object affordance	Press button to open door	93.3 ± 14.9

### Evolution

