USCUniversity of Southern California

- Prior methods train policies on large datasets or exploit primitive actions.
- key regions within the scene and reconstruct a voxel grid.



- Input: RGB-D images, language goals, and proprioception data.
- VLMs output pose of the object. We use this information to determine the language
- and arm ID.



VoxAct-B: Voxel-Based Acting and Stabilizing Policy for Bimanual Manipulation

I-Chun Arthur Liu, Sicheng He, Daniel Seita*, Gaurav S. Sukhatme*

Notivation

Bimanual manipulation tasks are challenging due to high dimensional action spaces.

How can we design a sample-efficient method without relying on primitive actions? • A voxel-based method that leverages Vision Language Models (VLMs) to prioritize

goal, roles of each arm (i.e., acting or stabilizing), and where to reconstruct a voxel grid. • **Output:** Discretized pose of the next best voxel, gripper action, collision avoidance flag,

Acting and Stabilizing Policies $\pi_a \pi_s$

- next best voxel with spatial equivariance properties, enabling more efficient learning from multi-modal demonstrations.
- Learning to predict arm ID allows policies to learn to map the appropriate acting or stabilizing actions to a given arm during training.

arm as acting. Methods are evaluated on 25 episodes of unseen test data. (1) Comparison with Baselines (2) Multi-Task Results (3) Ablations

	Open Jar		Open Drawer		Put Item in Drawer		Hand Over Item	
Method	10	100	10	100	10	100	10	100
Bimanual PerActs	8.0		36.8		5.6		0.0	
Diffusion Policy	4.8	21.6	4.8	5.6	2.4	4.8	0.0	0.0
VoxPoser	8.0	8.0	32.0	32.0	4.0	4.0	0.0	0.0
ACT w/Transformers	4.0	30.4	12.8	28.0	8.8	44.8	1.6	7.2
VoxAct-B (ours)	40.0	59.2	73.6	72.8	39.2	49.6	19.2	14.4

11	21
4	-)

	Open Jar	Open Drawer	Put Ite
ACT w/Transformers	2.7	12.0	14.7
VoxAct-B (ours)	21.3	62.7	17.3



Open Jar

Success: 5 out of 10 trials





Acting and stabilizing formulation exploits the discretized action space that predicts the

Open Metho Drawer m VoxAct-B w/o VLMs 19.2 wer VoxAct-B w/o Segment Anything 67.2 VoxAct-B w/o acting and stabilizing 64.8 68.0 VoxAct-B w/o arm ID 73.6 VoxAct-B (ours)

Real-World Results





