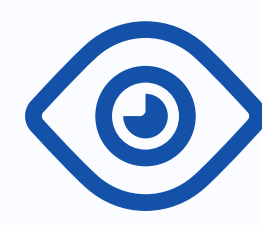


Baijun Chen^{5,2*}, Weijie Wan^{6*}, Tianxing Chen^{4*}, Xianda Guo^{7,2*}, Congsheng Xu¹, Yuanyang Qi³, Haojie Zhang³, Longyan Wu⁸, Tianling Xu¹, Zixuan Li⁶, Yizhe Wu³, Rui Li³, Xiaokang Yang¹, Ping Luo⁴, Wei Sui^{2†}, and Yao Mu^{1†}

¹ScaleLab, Shanghai Jiao Tong University ²D-Robotics ³ViTai Robotics ⁴The University of Hong Kong ⁵Nanjing University ⁶Shenzhen University ⁷Wuhan University ⁸Fudan University

1 Why Viso-Tactile Manipulation



Vision suffers from occlusion and poor contact observability.



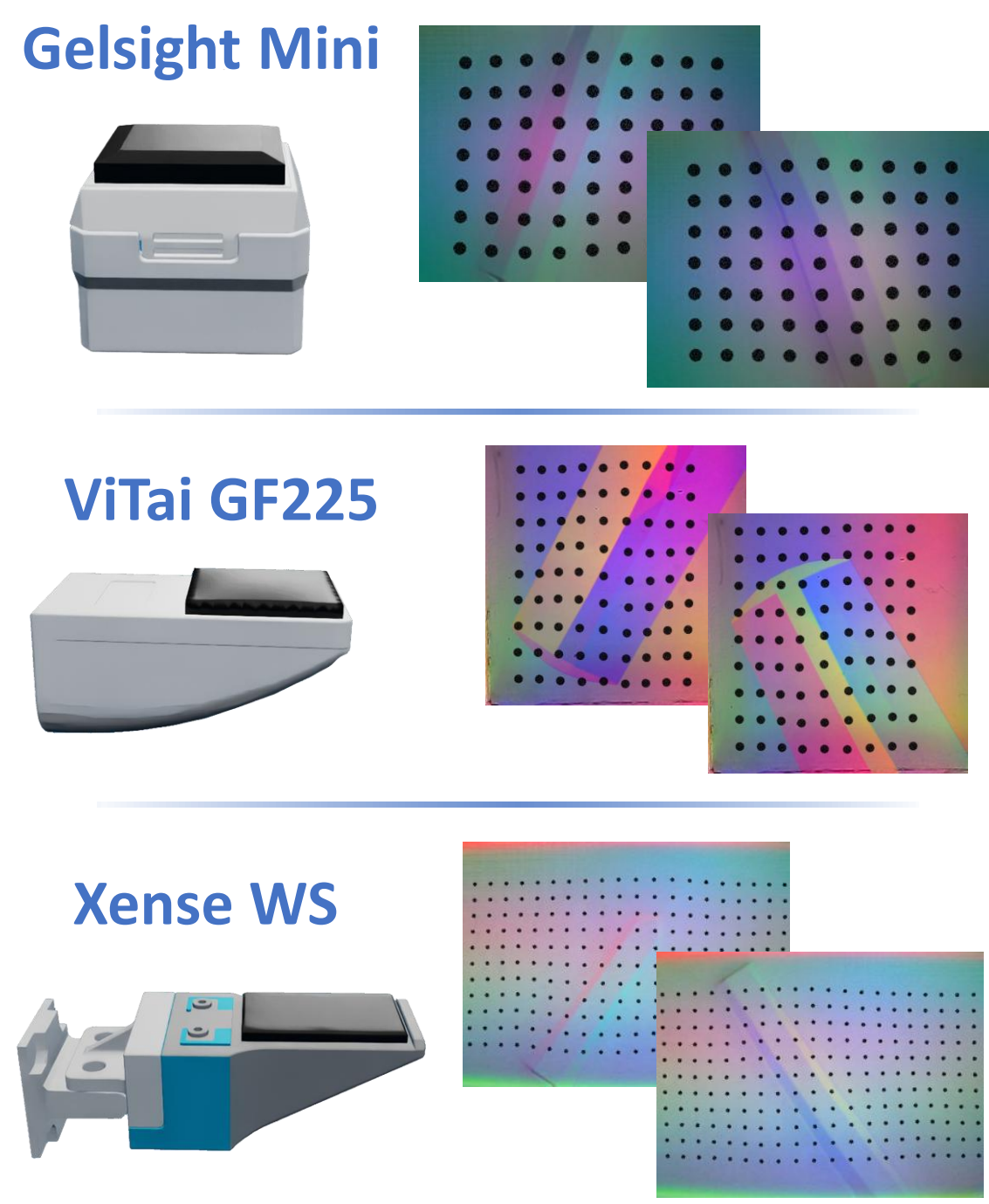
Tactile sensing provides local contact geometry, deformation, and shear cues.



Large-scale tactile data and standardized benchmarks are scarce.

2 Scalable Visuo-Tactile Simulation

3 Visuo-Tactile Sensors

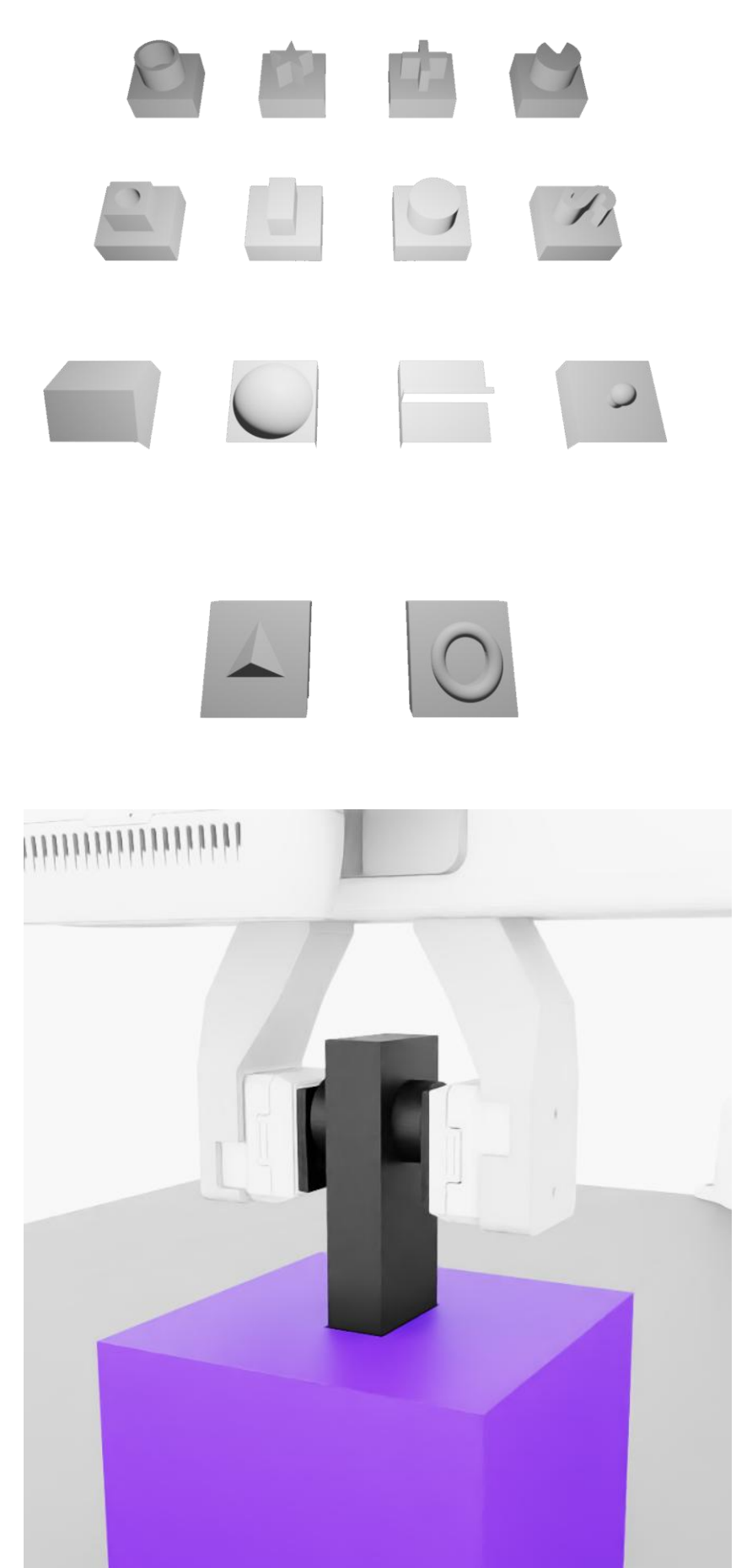


Tactile-Aware Manipulation Primitives

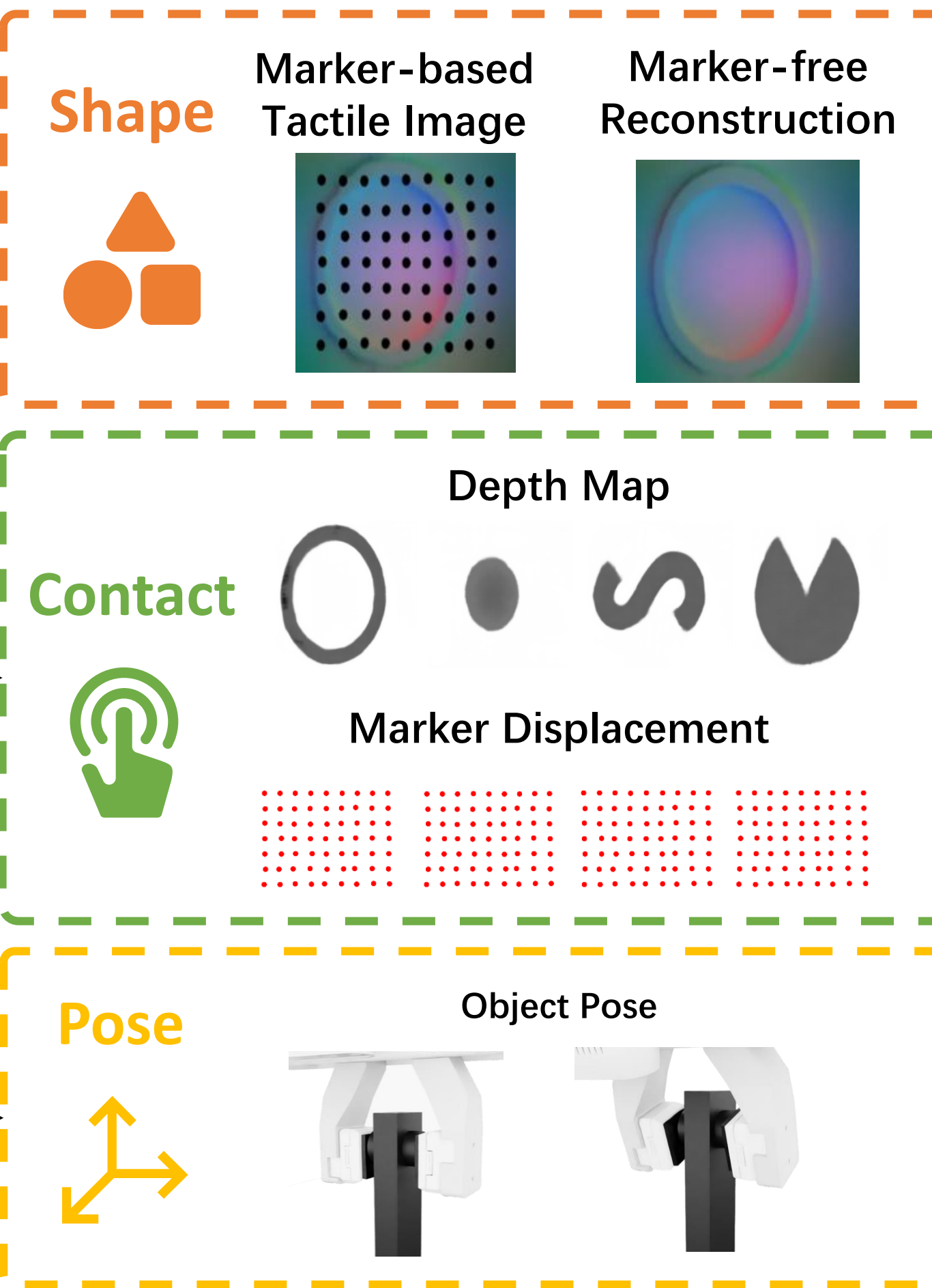
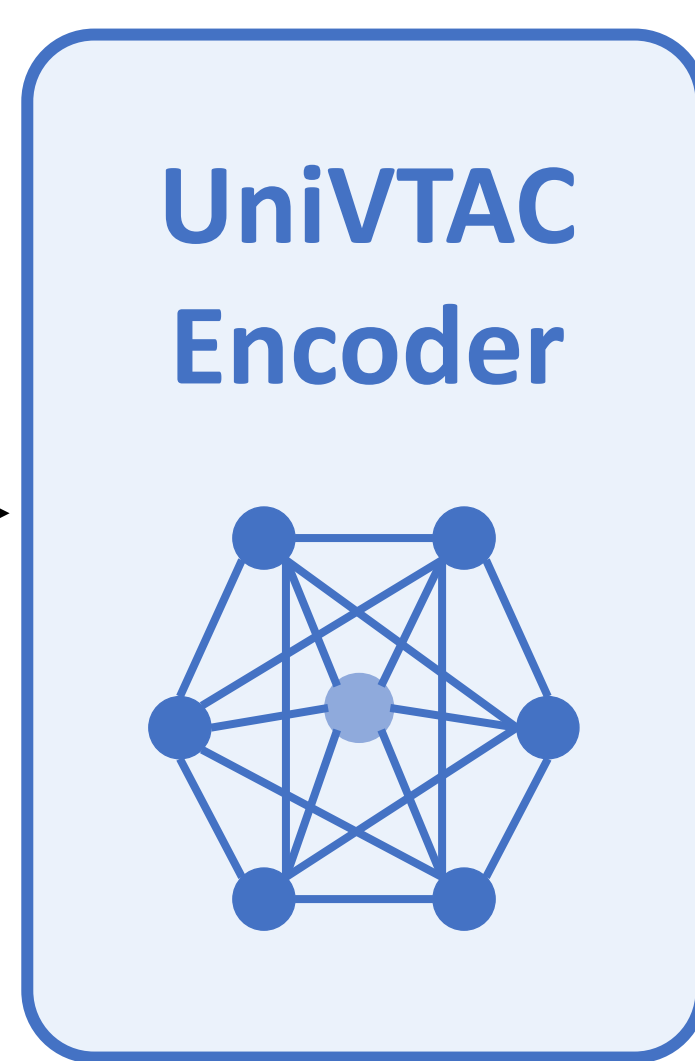
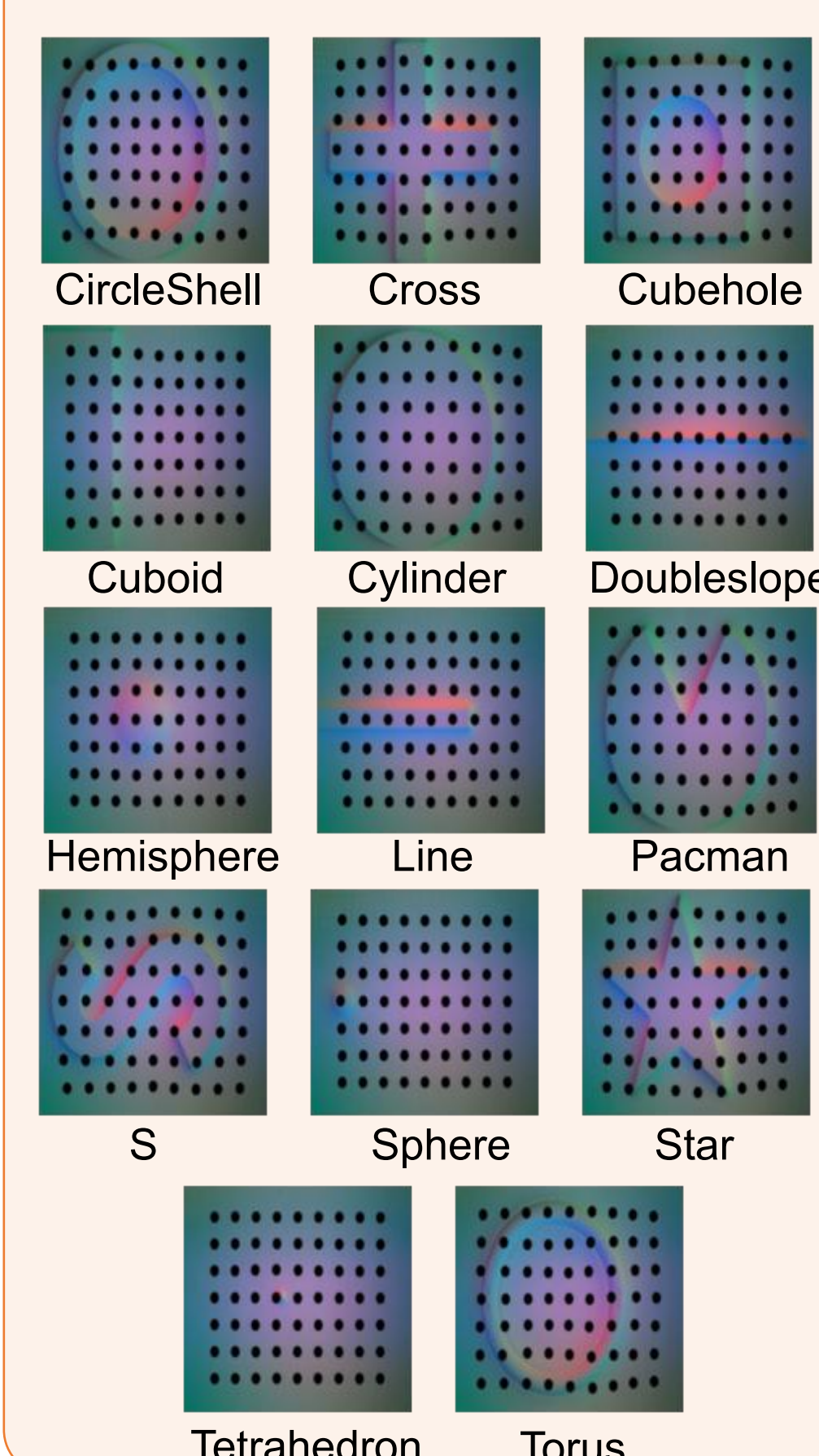


Unified simulation platform for visuo-tactile data generation across diverse sensors.

3 Privileged Supervision & UniVTAC Encoder



Raw Visuo-Tactile Observation

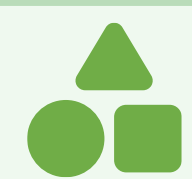
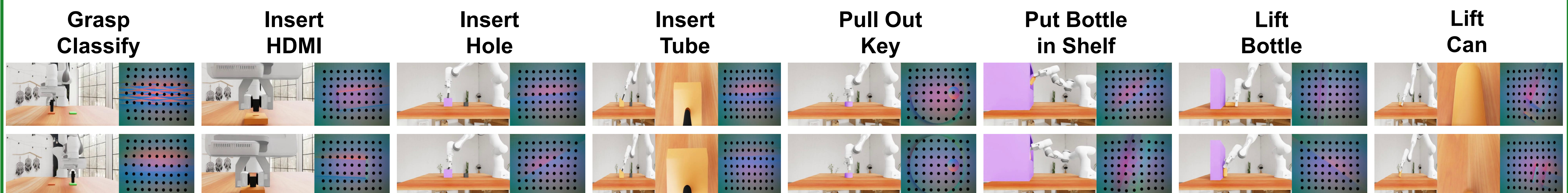


Simulation provides privileged physical annotation that are difficult to obtain in the real world.



Pretrain in simulation, then integrate the encoder into downstream policy learning.

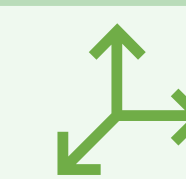
4 UniVTAC Benchmark



Shape preception



Pose reasoning



Contact-rich Interaction

5 Main Results



Simulation (8 Tasks)

ACT (Vision-Only) 30.9% → ViTAL(Strong Baseline) 40.5% → ACT+UniVTAC(Ours) 48.0% +17.1% ↑

Method	Lift Bottle	Pull-out Key	Lift Can	Put Bottle in Shelf	Insert Hole	Insert HDMI	Insert Tube	Grasp Classify	Average
ACT	42	28	20	28	19	15	45	50	30.9
ViTAL	72	47	8	32	25	6	34	100	40.5
Ours	71	46	29	31	24	28	56	99	48



Real-world (3 Tasks)

ACT (Vision-Only) 43.3% → ACT+UniVTAC (Ours) 68.3% +25.0% ↑

Method	Insert Tube	Insert USB	Bottle Upright	Average
ACT	55	15	60	43.3
ACT+UniVTAC	85	25	95	68.3

