

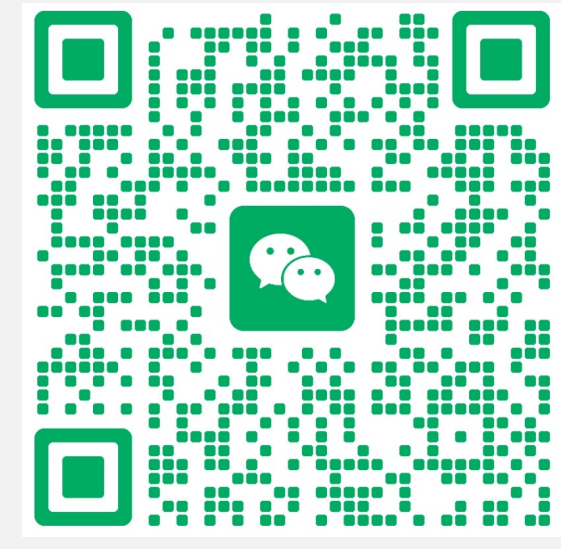


SpatialActor: Exploring Disentangled Spatial Representations for Robust Robotic Manipulation



清华大学
Tsinghua University

AAAI 26 Oral



Hao Shi
(石昊)
@THU

Hao Shi¹, Bin Xie², Yingfei Liu², Yang Yue¹, Tiancai Wang²,
Haoqiang Fan², Xiangyu Zhang^{3,4}, Gao Huang¹✉

¹LeapLab, Tsinghua University, ²Dexmal, ³MEGVII, ⁴StepFun

✉ shi-h23@mails.tsinghua.edu.cn



Page

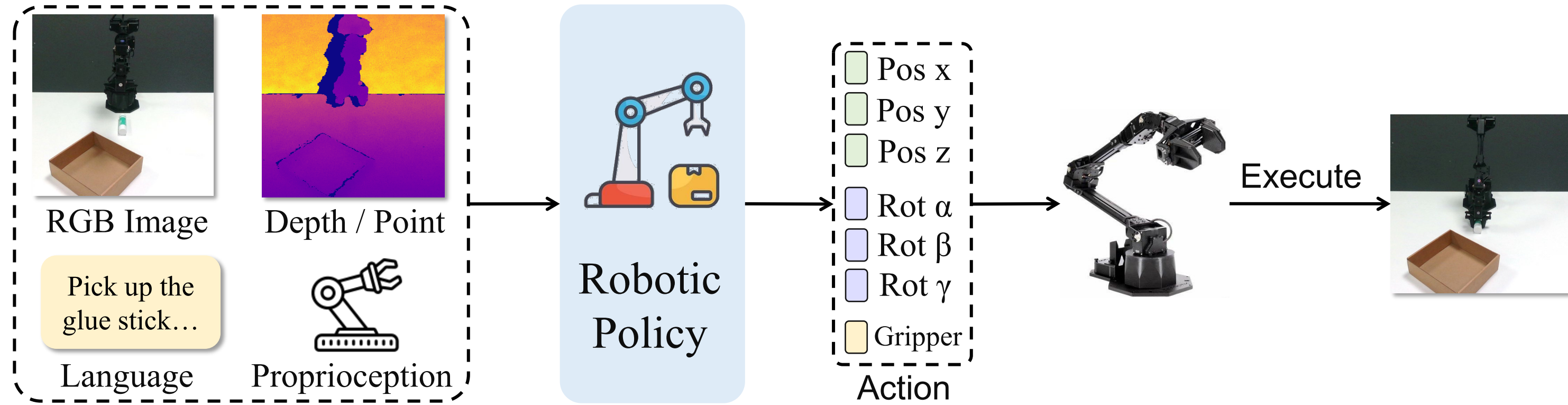


Code

<https://shihao1895.github.io>

1 Introduction

① What is Robotic Manipulation Task?



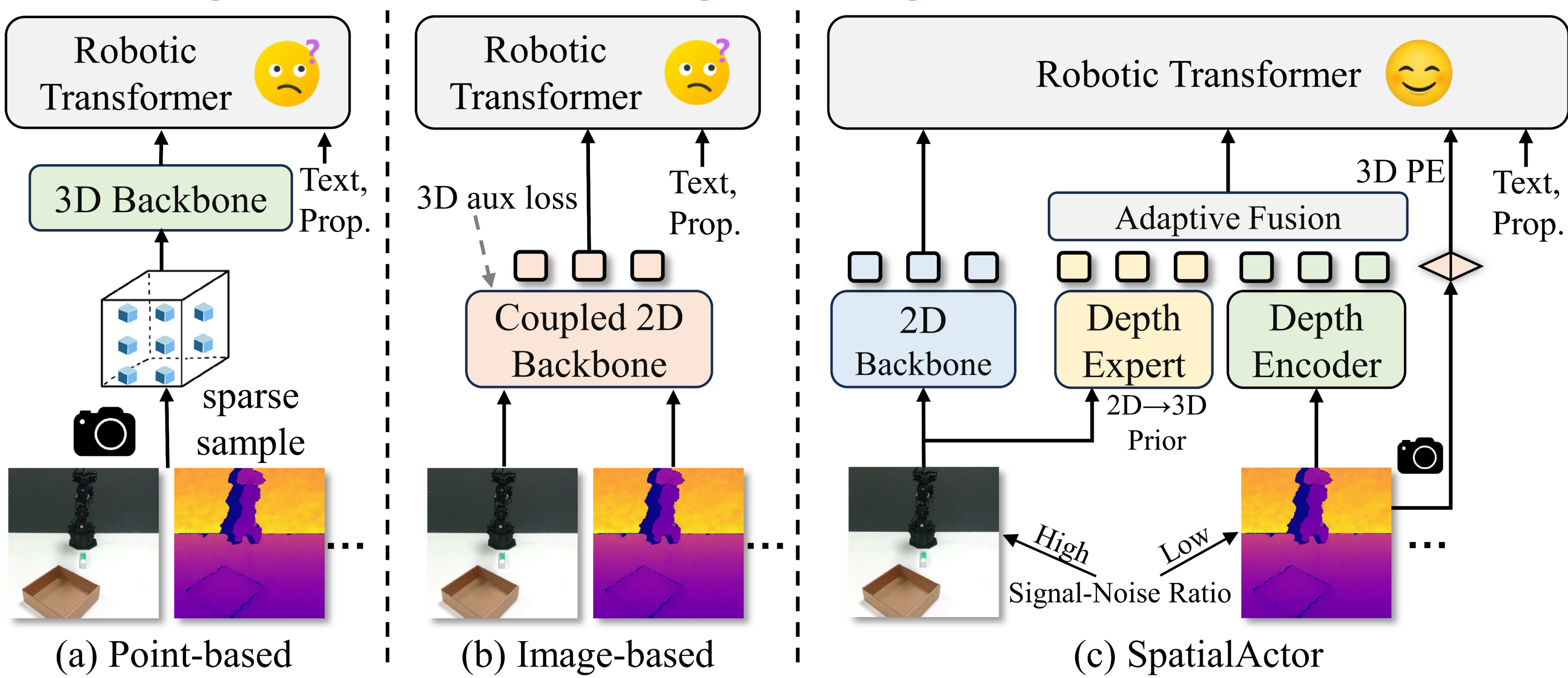
② Some Challenges in Robotic Manipulation

- ✓ Fine-grained spatial semantic understanding.
- ✓ Robustness to sensor noise.
- ✓ Low-level spatial structure inductive bias.



How to build robust spatial representations for robotic manipulation?

③ Comparison of Robotic Spatial Representations Solution



✓ Point-based

Sparse fusion → Loss fine-grained semantics.

✓ Image-based (e.g., RVT)

Coupled encode → Noisy depth interfere semantics.

✓ Ours: Disentangled framework

① Visual semantics

② Complementary high-level geometry

■ Fine-grained yet noisy: from real sensor depth

■ Robust but coarse-grained: from depth expert

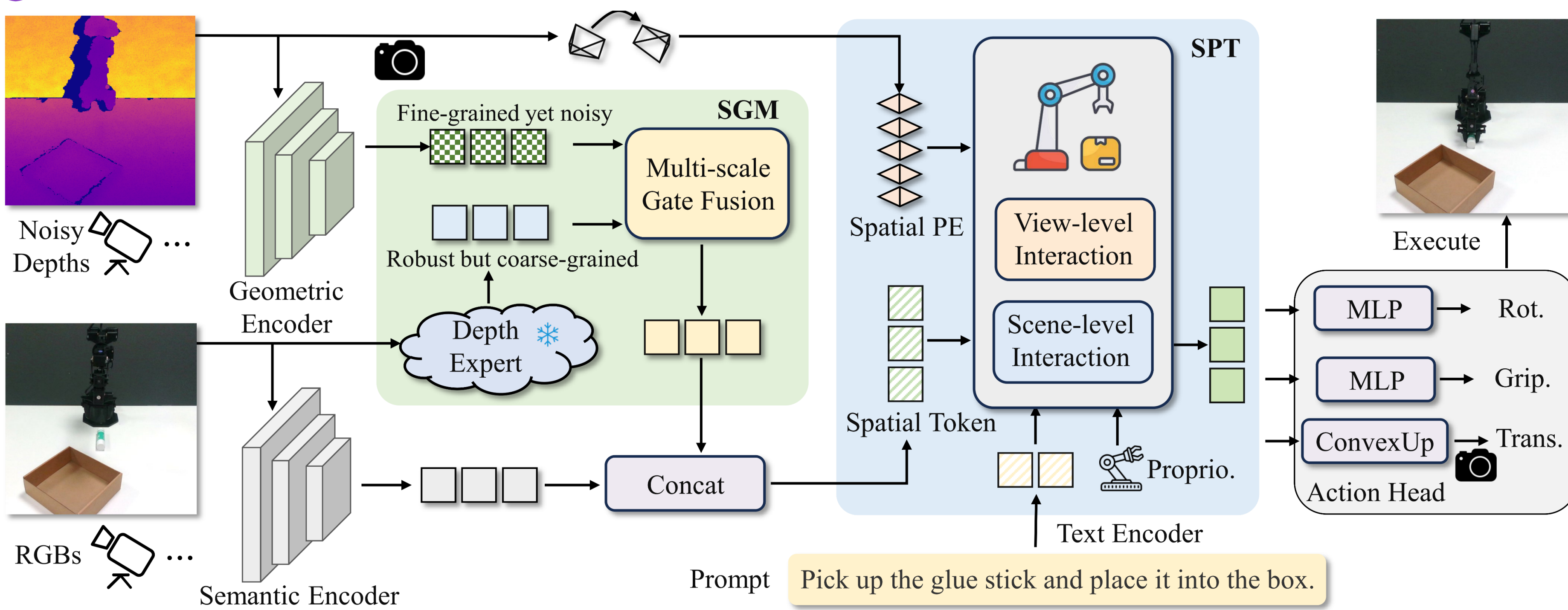
(High-SNR semantics + 2D-to-3D priors help smooth noisy depth)

③ Low-level geometry: explicit 2D-3D correspondence

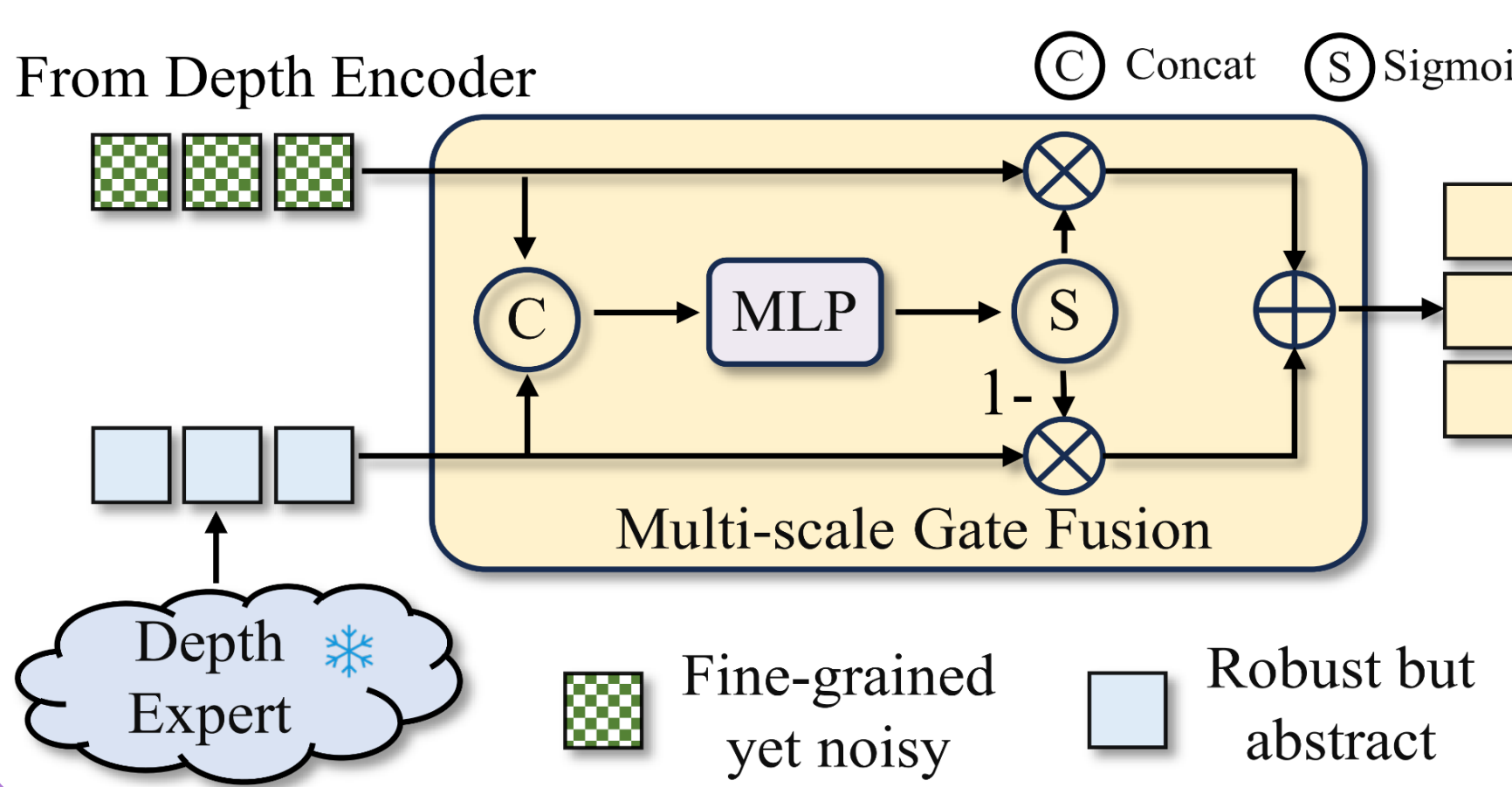
2 Method

3 Experiments

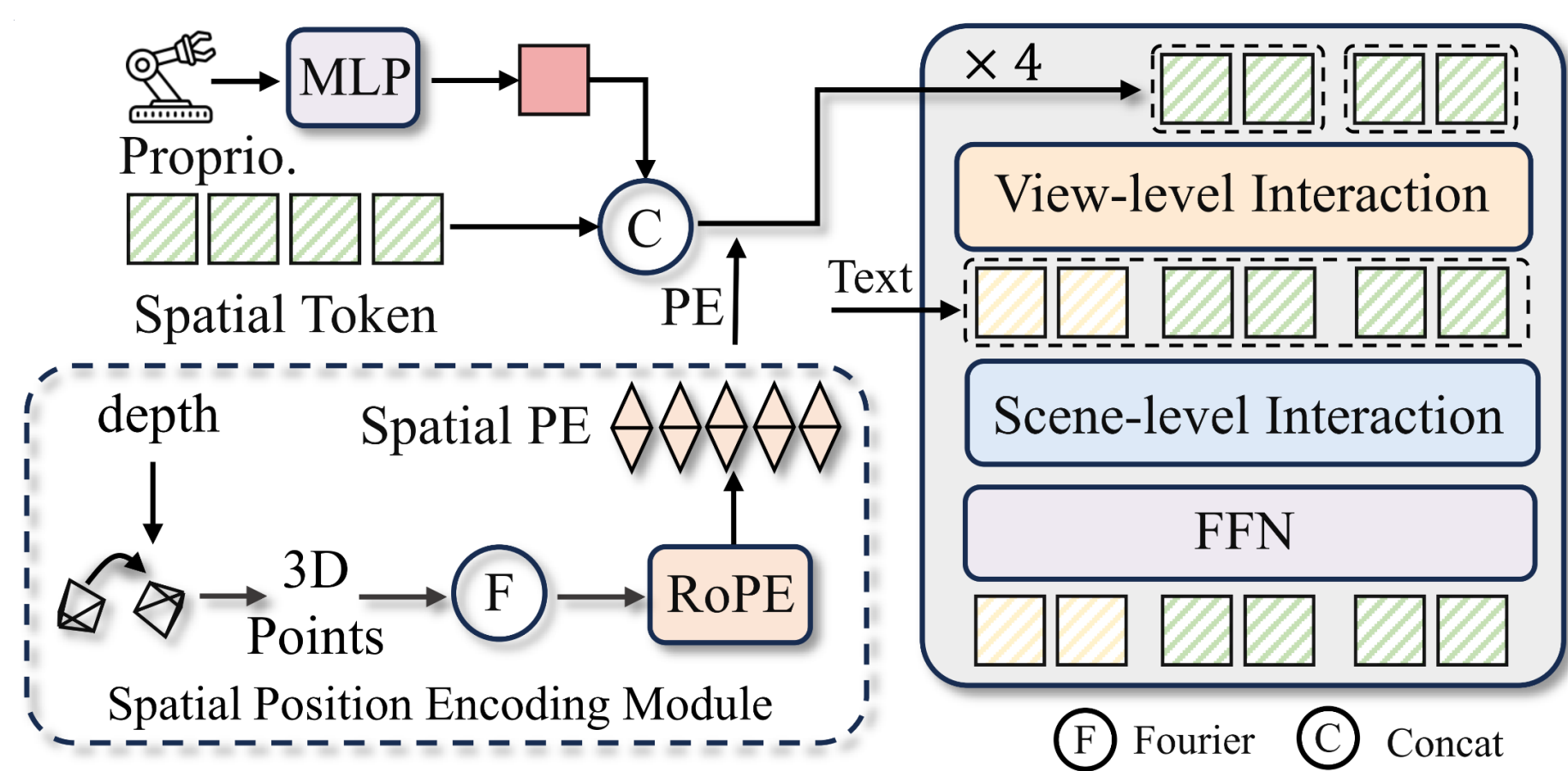
① Overall Framework



② Semantic-Guided Geometric Module



③ Spatial Transformer

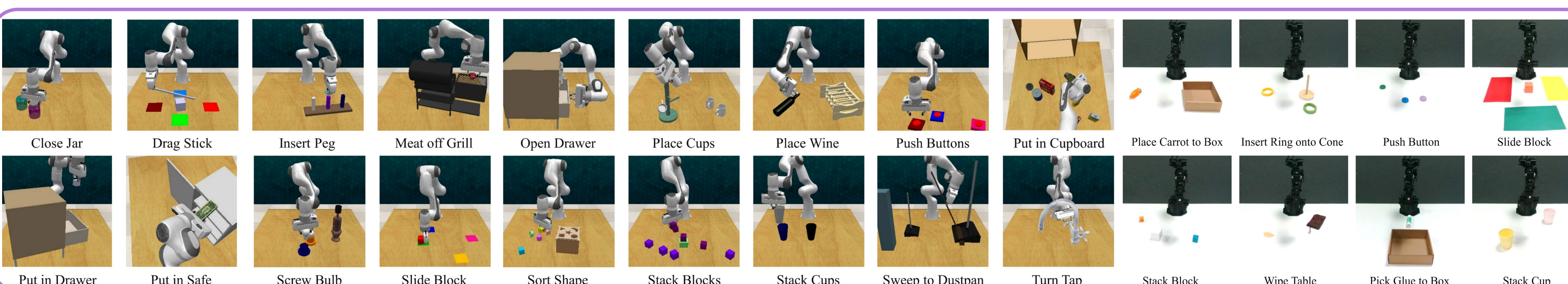


① RL Bench SOTA: 87.4% Success Rate

Models	Avg. Success ↑	Avg. Rank ↓	Close Jar	Drag Stick	Insert Peg	Meat off Grill	Open Drawer	Place Cups	Place Wine	Push Buttons
C2F-ARM-BC	20.1	9.5	24.0	24.0	4.0	20.0	20.0	0.0	8.0	72.0
HiveFormer	45.3	7.8	52.0	76.0	0.0	100.0	52.0	0.0	80.0	84.0
PolarNet	46.4	7.3	36.0	92.0	4.0	100.0	84.0	0.0	40.0	96.0
PerAct	49.4	7.1	55.2±4.7	89.6±4.1	5.6±4.1	70.4±2.0	88.0±5.7	2.4±3.2	44.8±7.8	92.8±3.0
RVT	62.9	5.3	52.0±2.5	99.2±1.6	11.2±3.0	88.0±2.5	71.2±6.9	4.0±2.5	91.0±5.2	100.0±0.0
Act3D	65.0	5.3	92.0	92.0	27.0	94.0	93.0	3.0	80.0	99.0
SAM-E	70.6	2.9	82.4±3.6	100.0±0.0	18.4±4.6	95.2±3.3	95.2±5.2	0.0±0.0	94.4±4.6	100.0±0.0
3D-Diff-Actor	81.3	2.8	96.0±2.5	100.0±0.0	65.0±4.1	96.8±1.6	89.6±4.1	24.0±7.6	93.6±4.8	98.4±2.0
RVT-2	81.4	2.8	100.0±0.0	99.0±1.7	40.0±0.0	99.0±1.7	74.0±11.8	38.0±4.5	95.0±3.3	100.0±0.0
SpatialActor	87.4±0.8	2.3	94.0±4.2	100.0±0.0	93.3±4.8	98.7±2.1	82.0±3.3	56.7±8.5	94.7±4.8	100.0±0.0

Models	Put in Cupboard	Put in Drawer	Put in Safe	Screw Bulb	Slide Block	Sort Shape	Stack Blocks	Stack Cups	Sweep Dustpan	Turn Tap
C2F-ARM-BC	0.0	4.0	12.0	8.0	16.0	8.0	0.0	0.0	0.0	68.0
HiveFormer	32.0	68.0	76.0	8.0	64.0	8.0	0.0	0.0	28.0	80.0
PolarNet	12.0	32.0	84.0	44.0	56.0	12.0	4.0	8.0	52.0	80.0
PerAct	28.0±4.4	51.2±4.7	84.0±3.6	17.6±2.0	74.0±3.0	16.8±4.7	26.4±3.2	2.4±2.0	52.0±0.0	88.0±4.4
RVT	49.6±3.2	88.0±5.7	91.2±3.0	48.0±5.7	81.6±5.4	36.0±2.5	28.8±3.9	26.4±8.2	72.0±0.0	93.6±4.1
Act3D	51.0	90.0	95.0	47.0	93.0	8.0	12.0	9.0	92.0	94.0
SAM-E	64.0±2.8	92.0±5.7	95.2±3.3	78.4±3.6	95.2±1.8	34.4±6.1	26.4±4.6	0.0±0.0	100.0±0.0	100.0±0.0
3D-Diff-Actor	85.6±4.1	96.0±3.6	97.6±2.0	82.4±2.0	97.6±3.2	44.0±4.4	68.3±3.3	47.2±8.5	84.0±4.4	99.2±1.6
RVT-2	66.0±4.5	96.0±0.0	96.0±0.0	88.0±4.9	92.0±2.8	35.0±7.1	80.0±2.8	69.0±5.9	100.0±0.0	99.0±1.7
SpatialActor	72.0±3.6	98.7±3.3	96.7±3.9	88.7±3.9	91.3±6.9	73.3±6.5	56±7.6	81.3±4.1	100.0±0.0	95.3±3.0

4 Visualization & Further Work

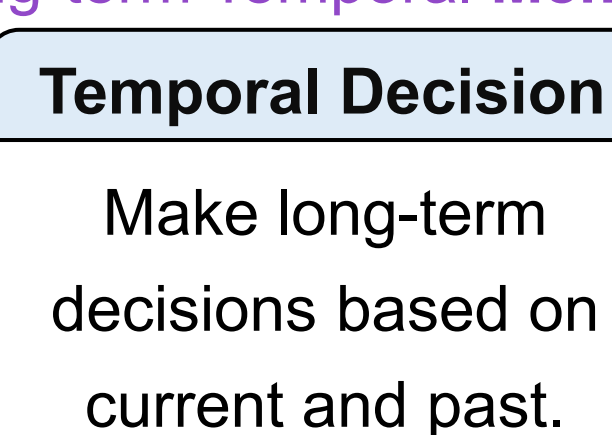
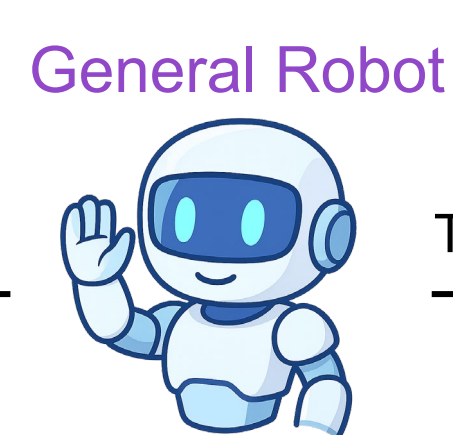
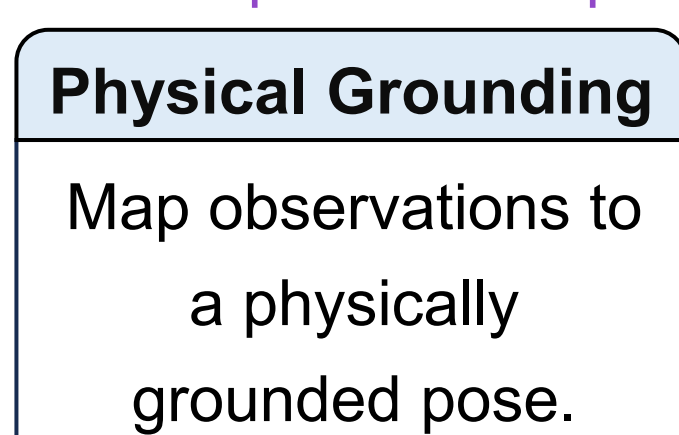


Equip Robot with Human-like Hippocampus Memory

Robust Spatial Perception

Long-term Temporal Memory

From Spatial to Temporal



MemoryVLA

④ Real-World Performance

