

# GMT: Goal-Conditioned Multimodal Transformer for 6-DOF Object Trajectory Synthesis in 3D Scenes

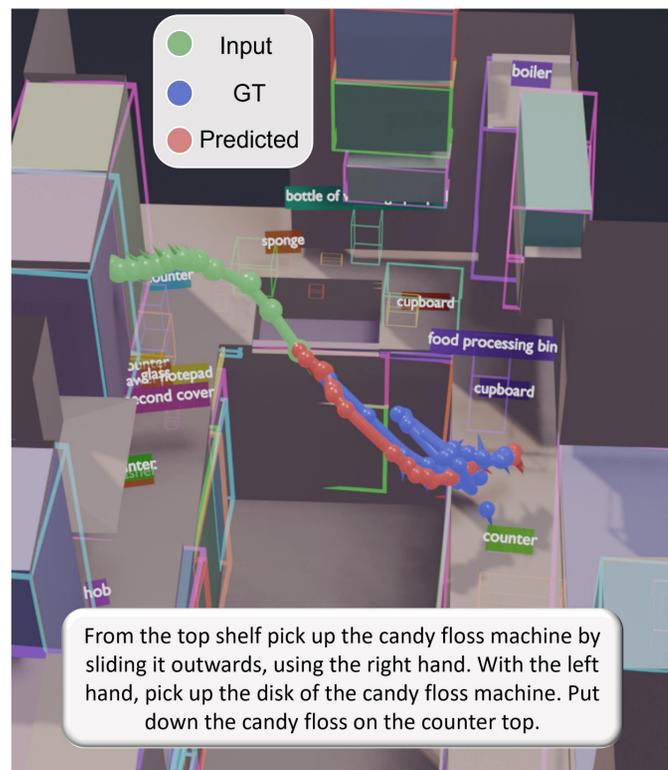
Huajian Zeng<sup>1,4,\*</sup>, Abhishek Saroha<sup>1,2,\*</sup>, Daniel Cremers<sup>1,2</sup>, Xi Wang<sup>1,2,3</sup>

<sup>1</sup>Technical University of Munich (TUM), <sup>2</sup>Munich Center for Machine Learning (MCML), <sup>3</sup>ETH Zürich, <sup>4</sup>Mohamed bin Zayed University of Artificial Intelligence (MBZUAI)

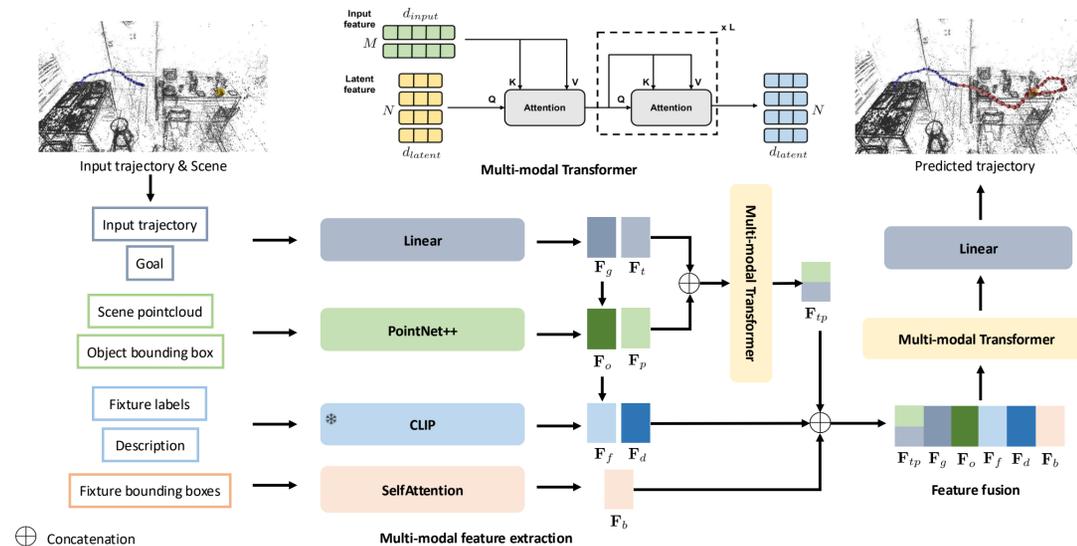


## Why Object Trajectory

Synthesizing controllable 6-DOF object manipulation trajectories in 3D environments is essential for enabling robots to interact with complex scenes.



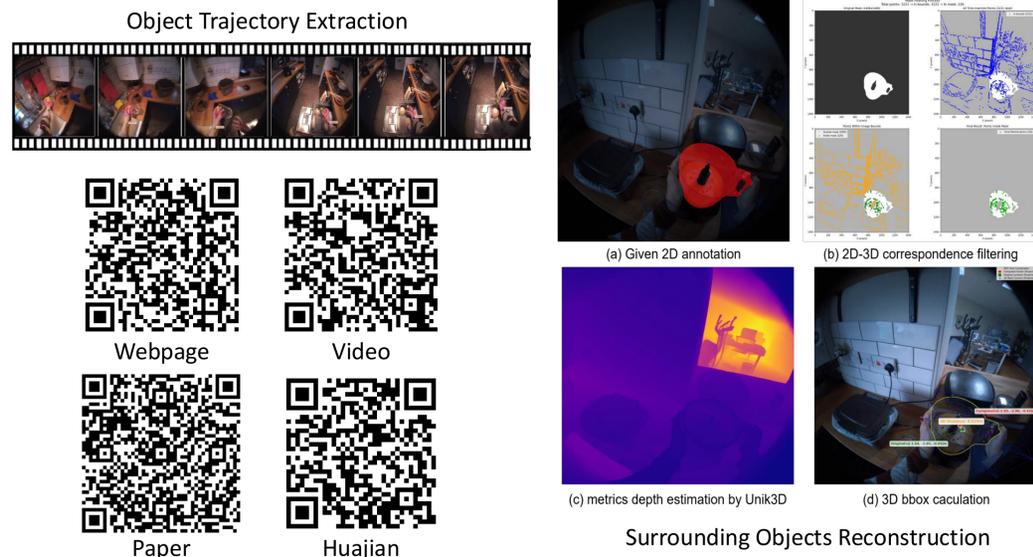
## GMT Architecture



## Key Design

**Object-centric generation:** predict object motion directly instead of modeling it as a byproduct of human motion.  
**Tailored conditioning:** combine local geometry, fixture semantics, language description, and goal pose.  
**Geometry-first fusion:** hierarchical transformer fusion improves spatial feasibility, long-horizon stability, and goal alignment.

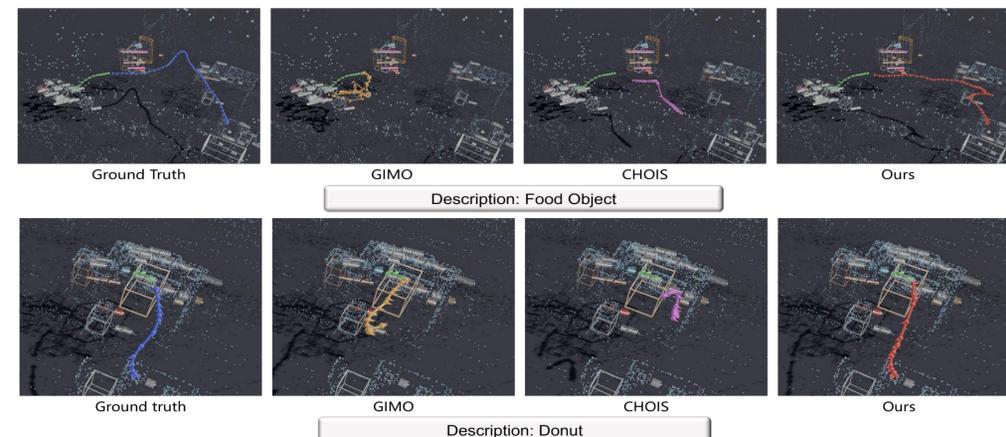
## Egocentric Video Preprocessing



## Experiments

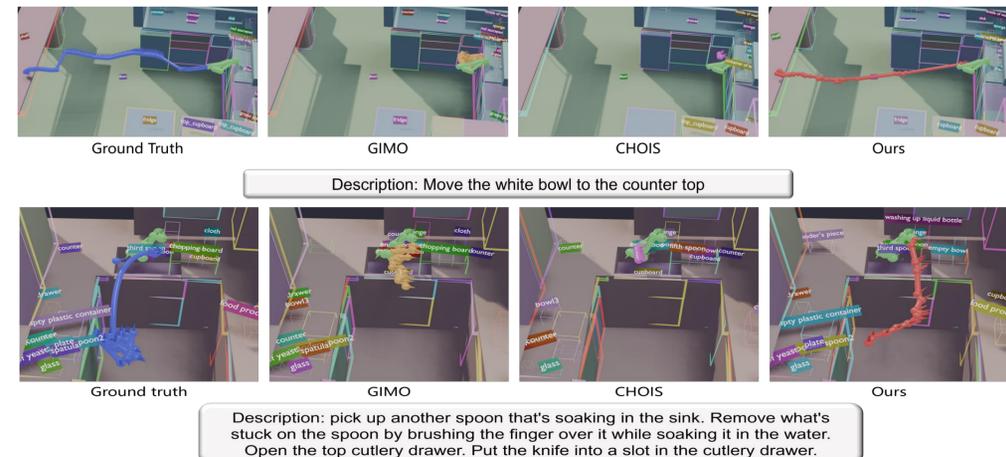
Controlled Scenarios (Aria Digital Twin Dataset [1])

Method	ADE[m] ↓	FDE[m] ↓	FD[m] ↓	AC[m] ↑	CR[%] ↓
GIMO [3]	0.982	1.401	1.511	0.140	19.6
CHOIS [4]	0.853	1.062	1.209	0.283	9.3
<b>GMT (Ours)</b>	<b>0.366</b>	<b>0.072</b>	<b>0.438</b>	<b>0.402</b>	13.1



Realistic Scenarios (HD-EPIC Dataset [2])

Method	ADE[m] ↓	FDE[m] ↓	FD[m] ↓	AC[m] ↑	CR[%] ↓
GIMO [3]	0.411	0.654	0.780	0.002	11.8%
CHOIS [4]	0.446	0.589	0.760	0.009	12.0%
<b>GMT (Ours)</b>	<b>0.283</b>	<b>0.034</b>	<b>0.391</b>	<b>0.037</b>	<b>10.3%</b>



## References

- [1] Pan X, Charron N, Yang Y, et al. Aria digital twin: A new benchmark dataset for egocentric 3d machine perception. CVPR 2023.
- [2] Perrett T, Darkhalil A, Sinha S, et al. Hd-epic: A highly-detailed egocentric video dataset. CVPR 2025.
- [3] Zheng Y, Yang Y, Mo K, et al. Gimo: Gaze-informed human motion prediction in context. ECCV 2022.
- [4] Li J, Clegg A, Mottaghi R, et al. Controllable human-object interaction synthesis. ECCV 2024.