



# SViMo: Synchronized Diffusion for Video and Motion Generation in Hand-object Interaction Scenarios

Lingwei Dang <sup>1\*</sup>, Ruizhi Shao <sup>2\*</sup>, Hongwen Zhang <sup>3</sup>, Wei MIN <sup>4</sup>, Yebin Liu <sup>2</sup>, Qingyao Wu <sup>†1</sup>

<sup>1</sup> South China University of Technology, <sup>2</sup> Tsinghua University, <sup>3</sup> Beijing Normal University, <sup>4</sup> Shadow AI

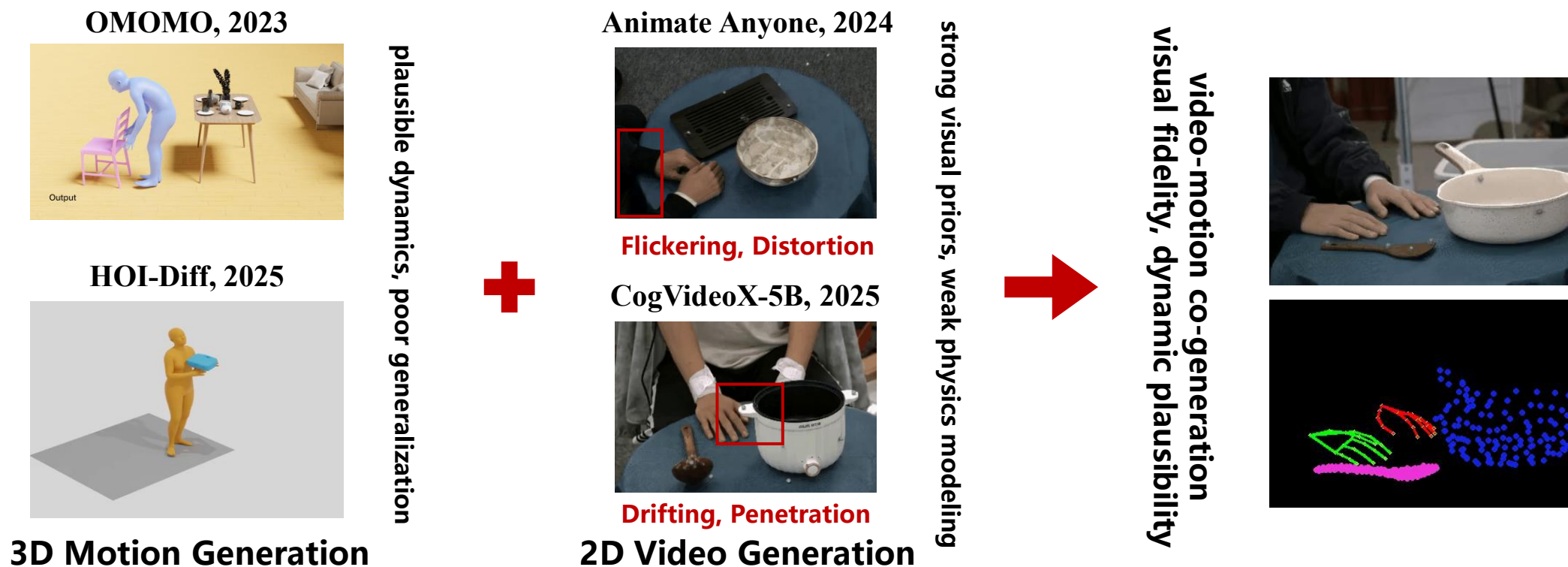
NeurIPS 2025 Spotlight

\* Equal contributions.

<sup>†</sup> Corresponding Author. Email: qyw@scut.edu.cn.

# 1. Background & Motivation

- 3D motion generation models produce plausible dynamics but suffer from limited data and **poor generalization**
- 2D Video generation models possess rich visual priors yet **lack awareness of physics**.

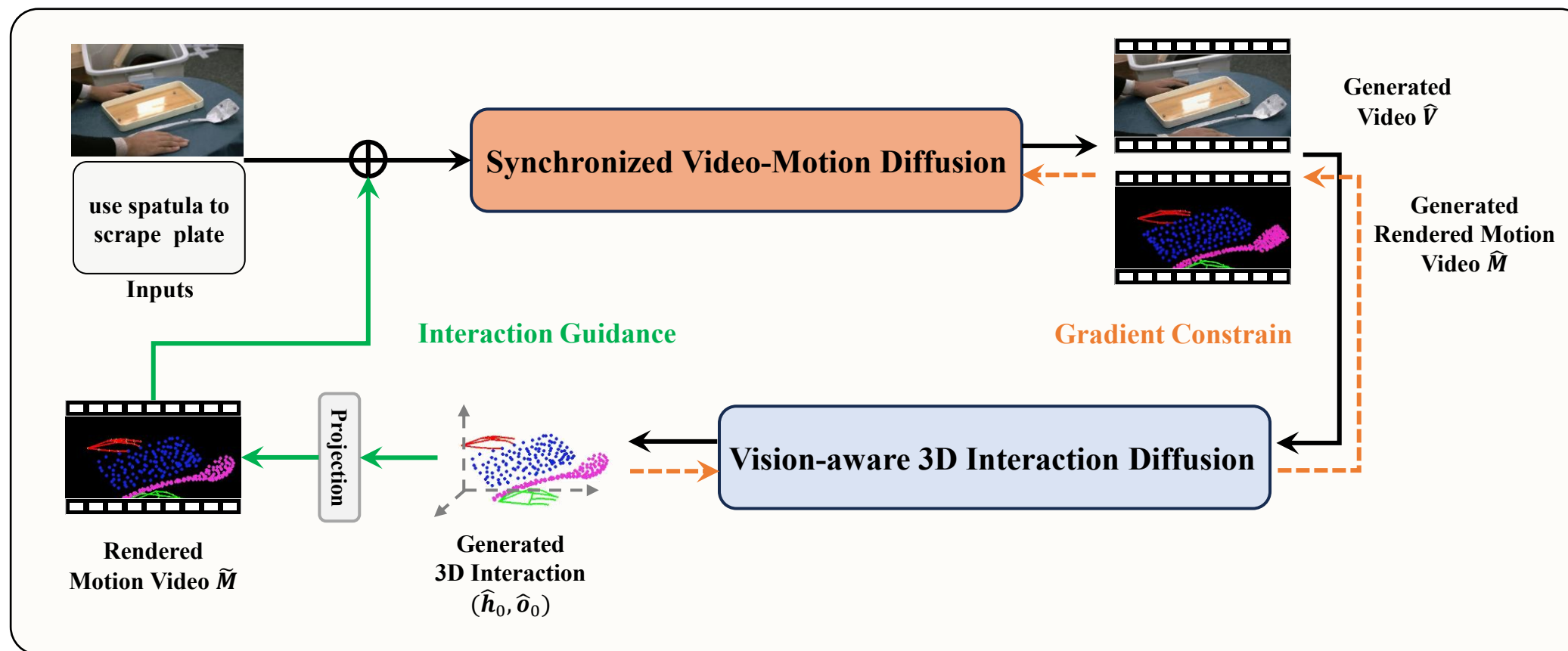


Object Motion Guided Human Motion Synthesis, SIGGRAPH Asia 2023; HOI-Diff: Text-Driven Synthesis of 3D Human-Object Interactions using Diffusion Models, CVPR 2025 Workshop; Animate Anyone: Consistent and Controllable Image-to-Video Synthesis for Character Animation, CVPR 2024; Cogvideox: Text-to-video diffusion models with an expert transformer, ICLR 2025.

**Visual appearance and motion dynamics share the same physical laws. We propose to unify visual priors and kinematic constraints through synchronized video-motion co-generation.**

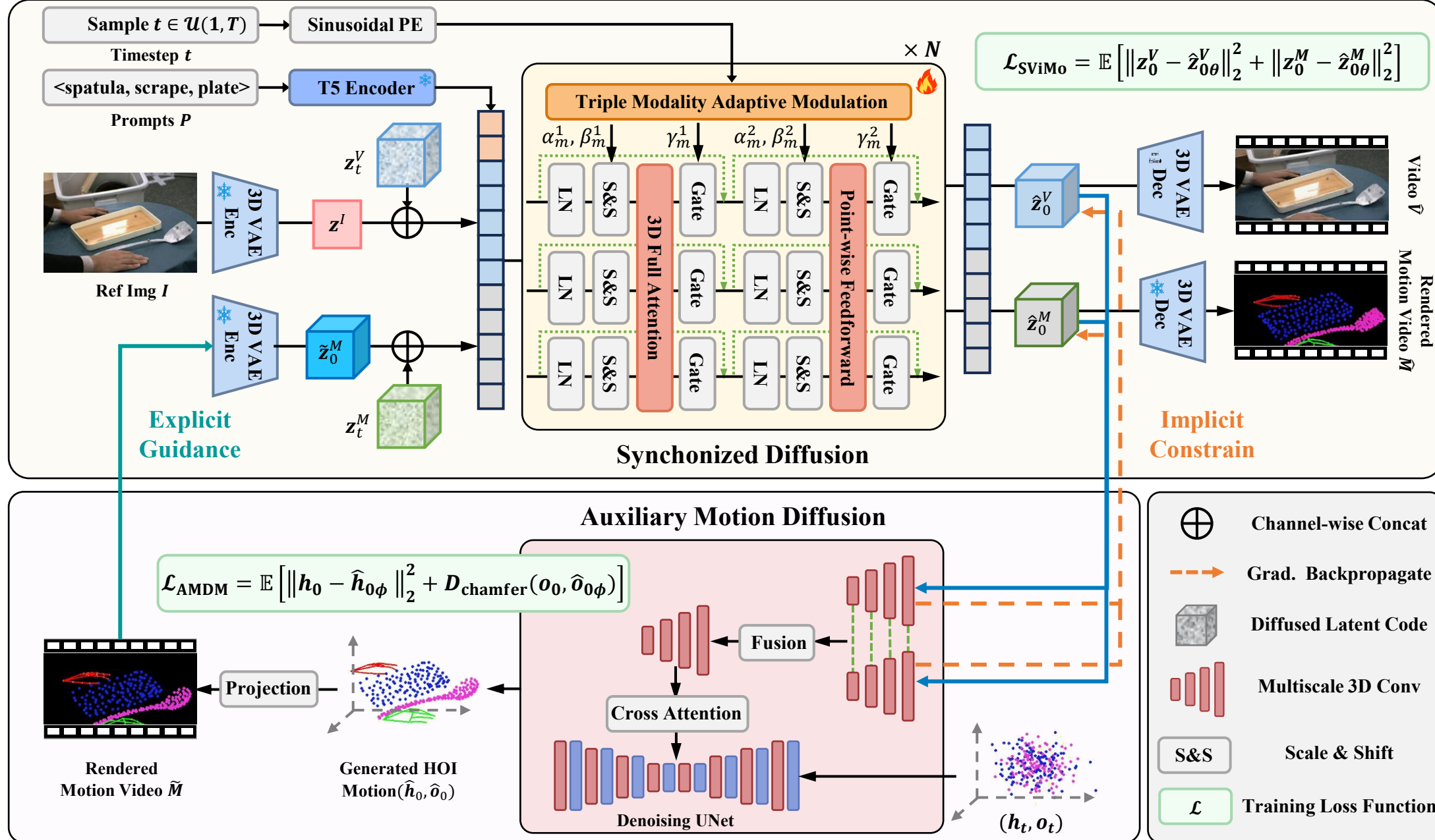
# 2. Method

## SViMo: A synchronized diffusion model for HOI video and motion joint generation



- ① End-to-end video-motion synthesis;      ② Visual realism and dynamic plausibility;      ③ Generalization ability

# 2. Method





# 3. Demonstrations of Our Method: Case 1

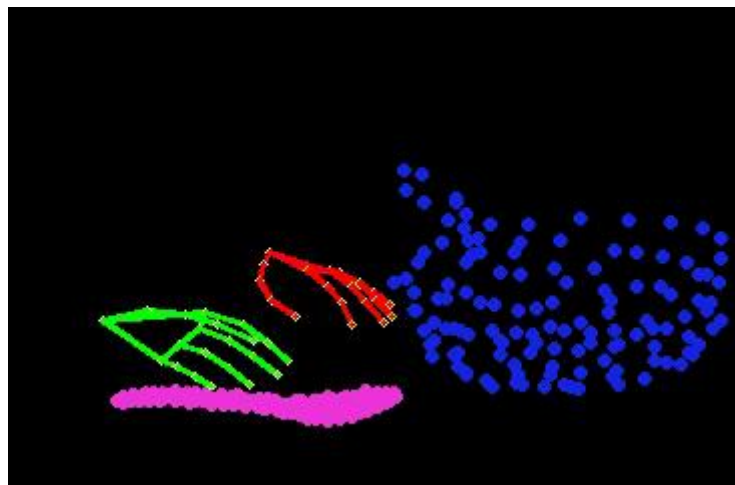
use spatula to scrape off pan

Generated Video



Realistic Video

Generated Motion



Plausible Motion

Overlaid Results



Consistent Video and Motion





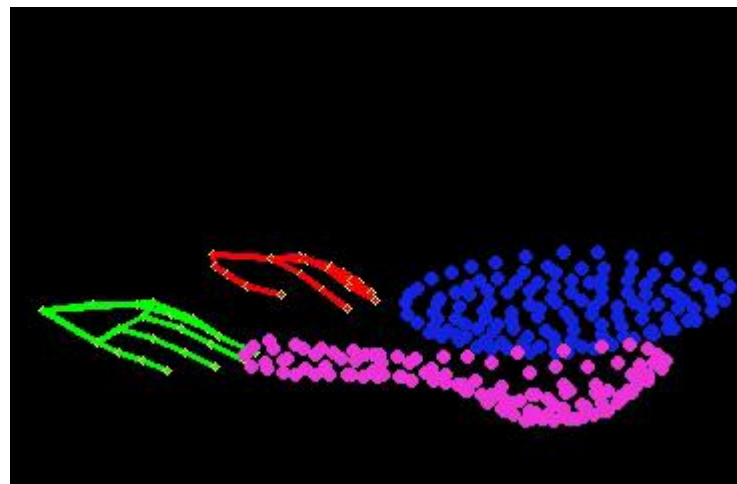
# 3. Demonstrations of Our Method: Case 2

use spoon to scrape off plate

Generated Video



Generated Motion



Overlaid Results



Realistic Video

Plausible Motion

Consistent Video and Motion



# 3. Demonstrations of Our Method: Case 3

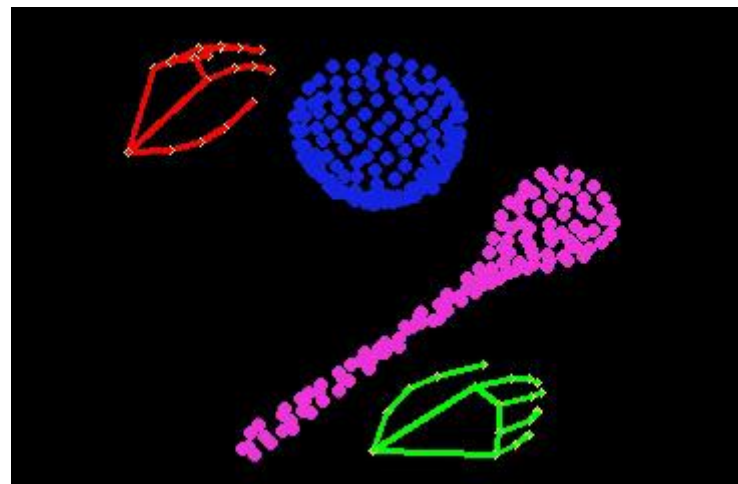
use spatula to put out bowl

Generated Video



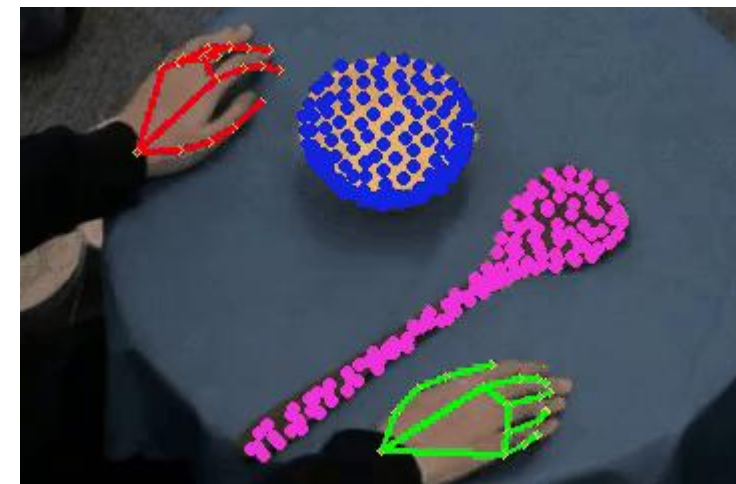
Realistic Video

Generated Motion



Plausible Motion

Overlaid Results



Consistent Video and Motion



# 4. Comparison of Videos: Case 1

use bowl to put in plate

Hunyuan-13B-Zeroshot



**Low-dynamic, Hallucination**

Animate Anyone



**Flickering, Distortion**

CogVideoX-5B



**Implausible movements**

Wan-14B-Zeroshot



**Hallucination, Implausible movements**

Easy Animate



**Flickering, Object inconsistency**

Ours







# 4. Comparison of Videos: Case 2

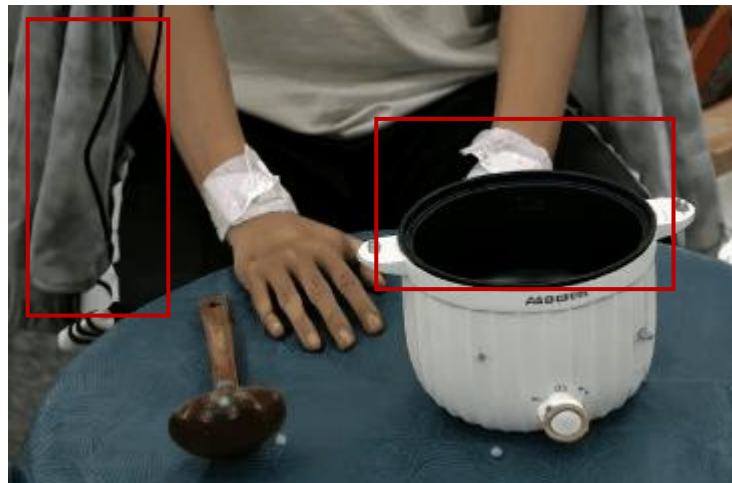
use spoon to scrape off pan

Hunyuan-13B-Zeroshot



Low-dynamic

Animate Anyone



Flickering, Distortion

CogVideoX-5B



Implausible movements, Penetration

Wan-14B-Zeroshot



Hallucination

Easy Animate



Flickering, Implausible movements

Ours

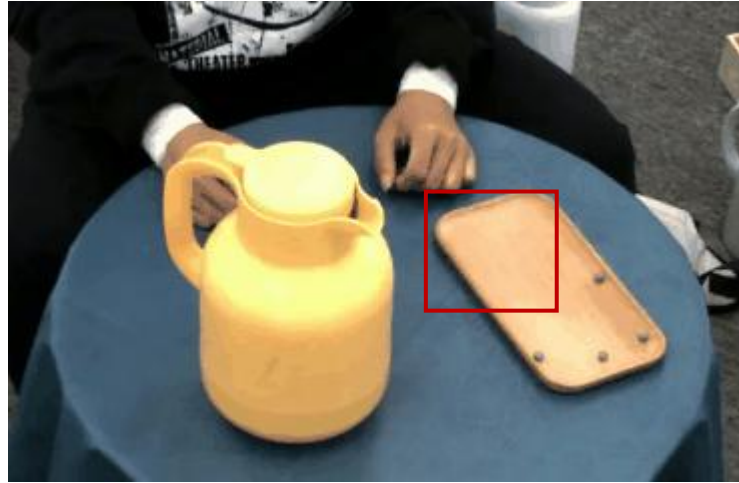




# 4. Comparison of Videos: Case 3

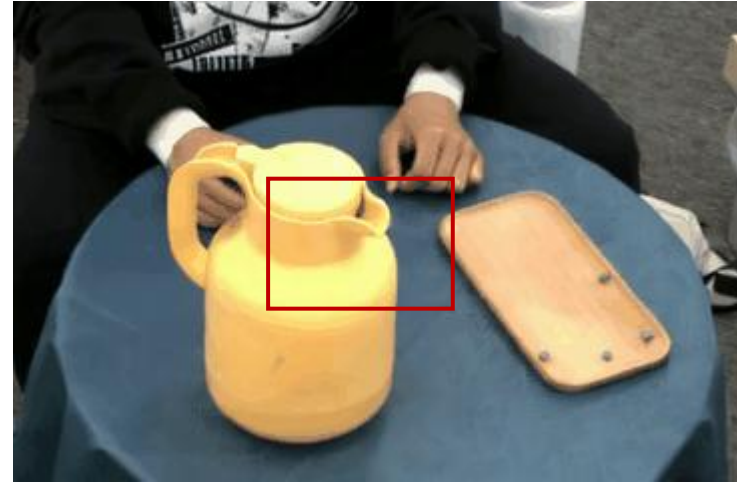
use kettle to pour in plate

Hunyuan-13B-Zeroshot



Low-dynamic, Hallucination

Animate Anyone



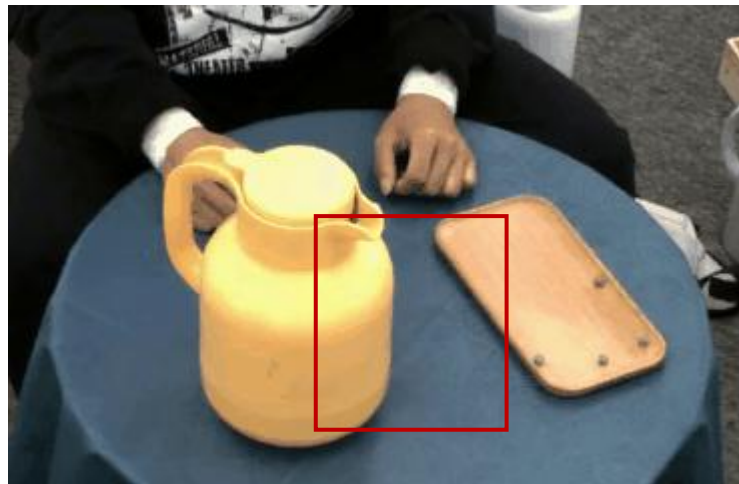
Flickering, Distortion

CogVideoX-5B



Object inconsistency

Wan-14B-Zeroshot



Hallucination, Camera shake

Easy Animate



Hallucination

Ours





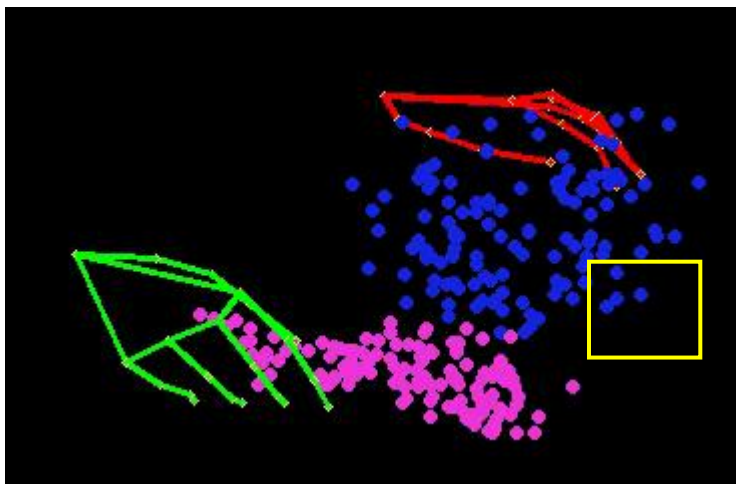


# 5. Comparison of Motions: Case 1

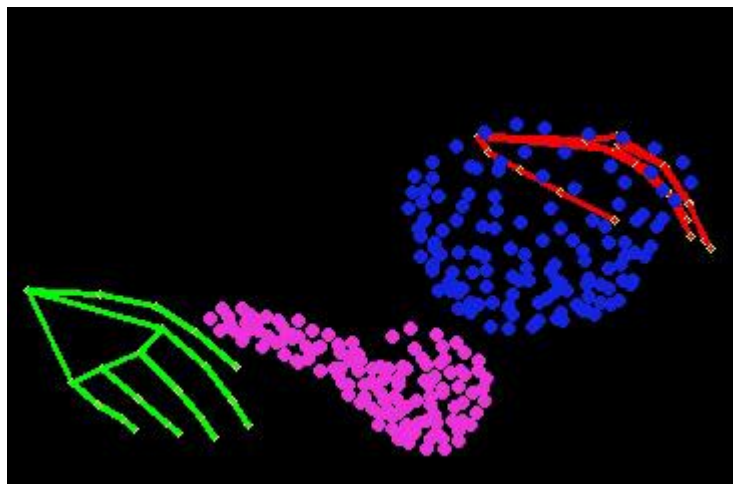
use roller to dust kettle

Generated  
Motion

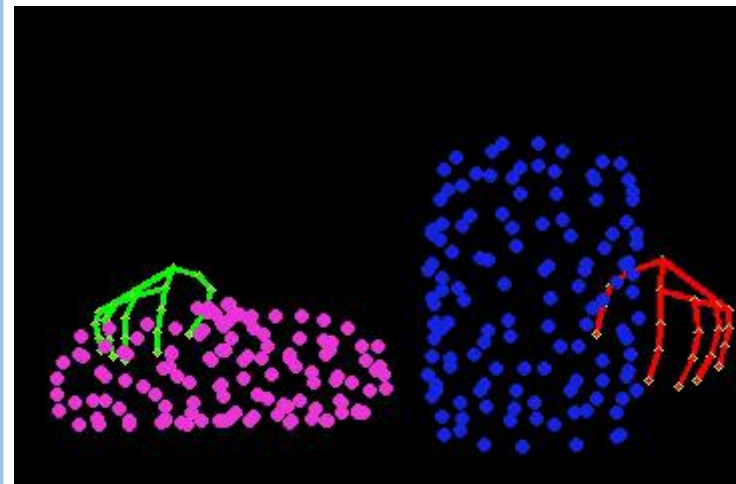
MDM



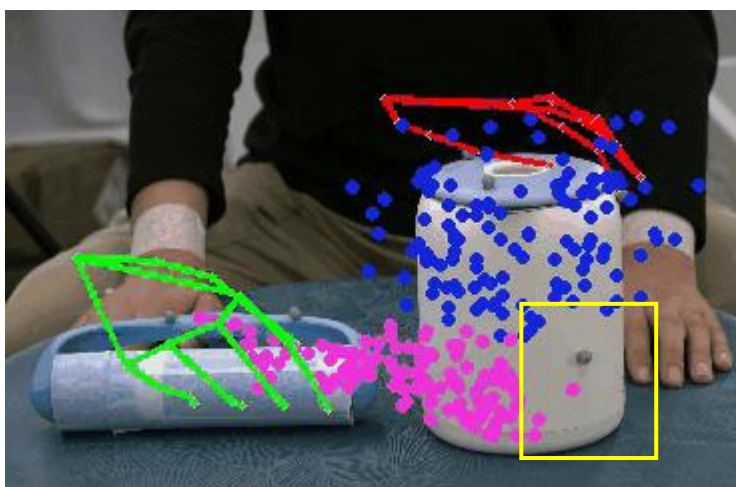
EMDM



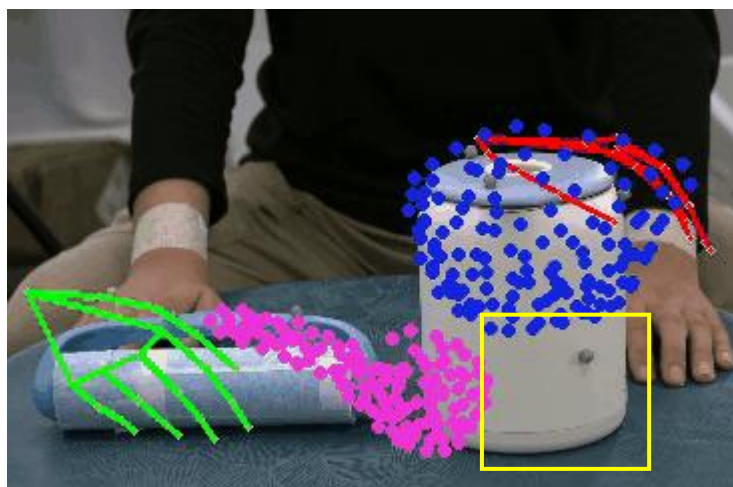
Ours



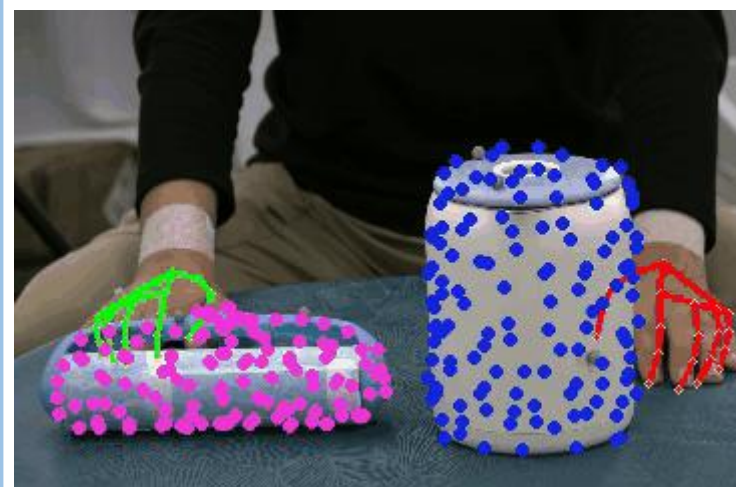
Overlaid  
Motion  
onto the Image



**Blurred outline, Low dynamic**



**Image-motion inconsistency**



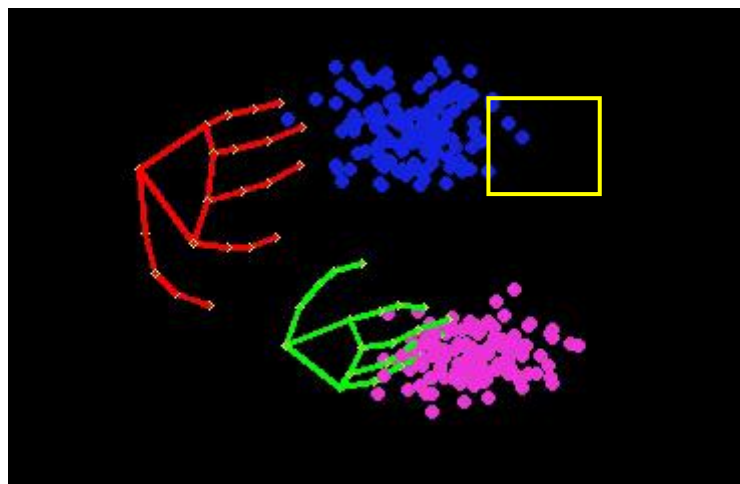


# 5. Comparison of Motions: Case 2

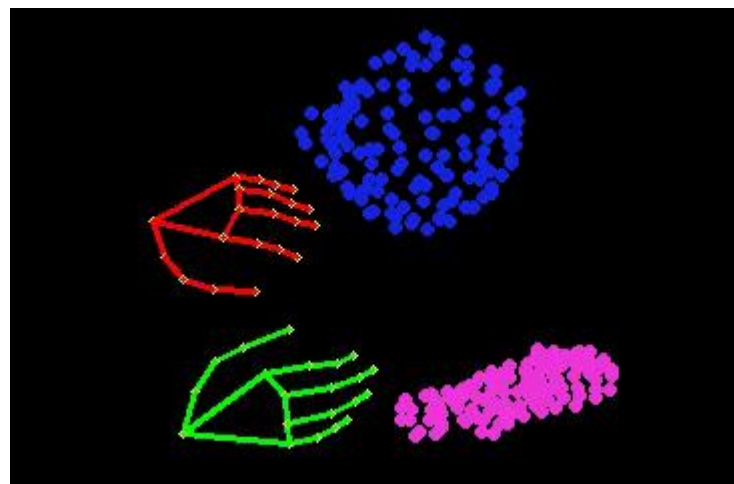
use spatula to stir plate

Generated  
Motion

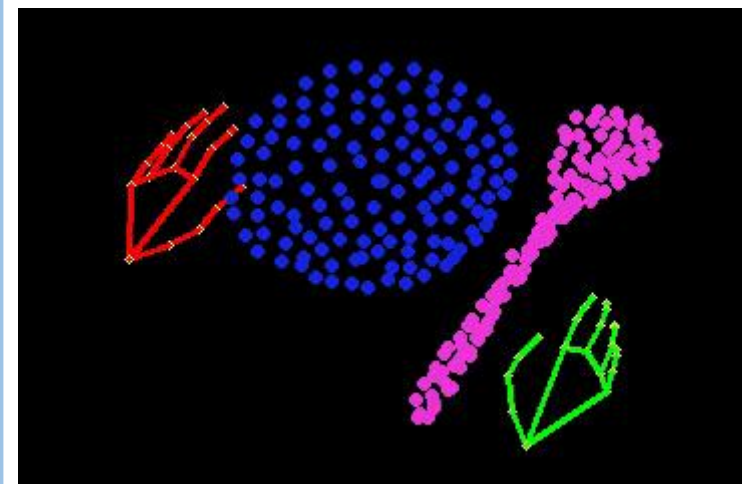
MDM



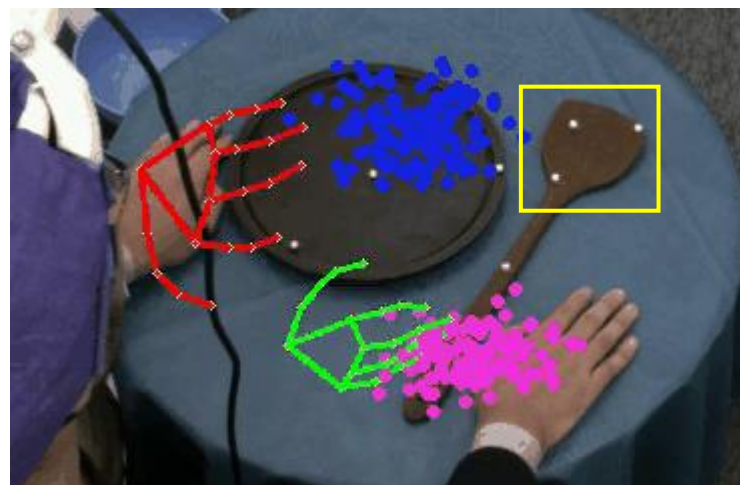
EMDM



Ours



Overlaid  
Motion  
onto the Image



Blurred outline, Implausible mov.

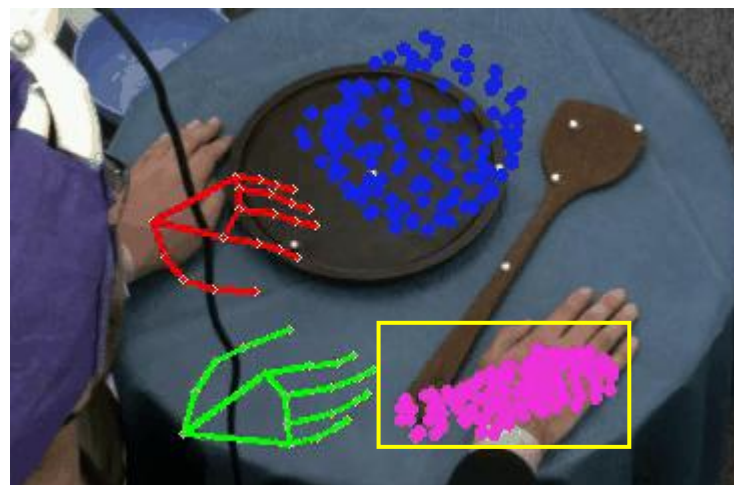
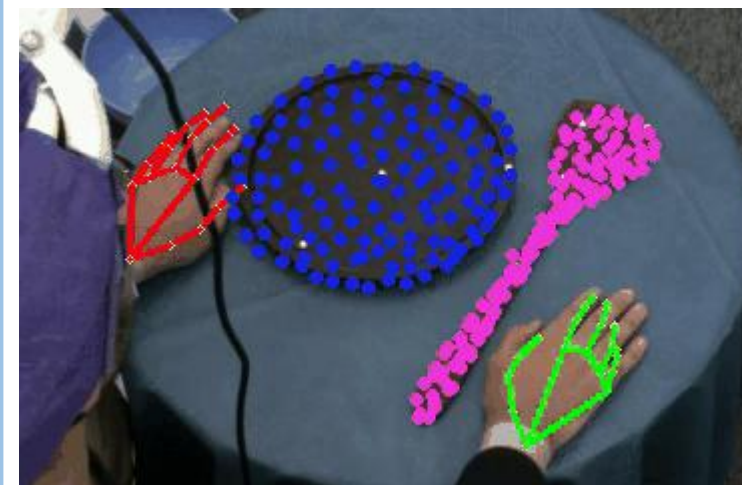


Image-motion inconsistency





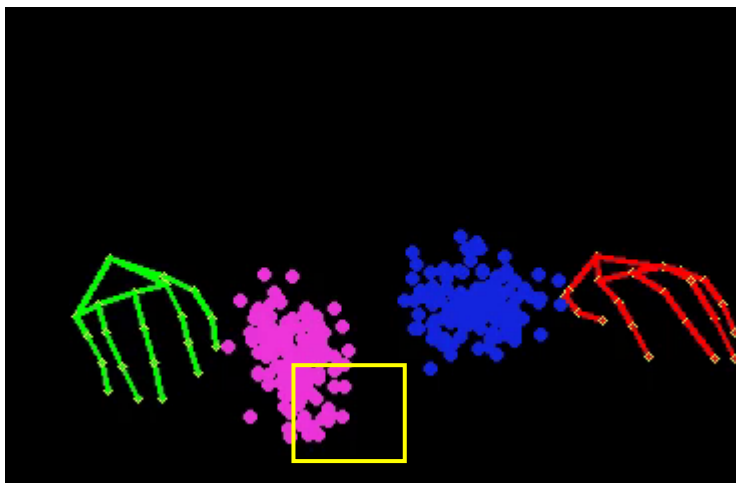


# 5. Comparison of Motions: Case 3

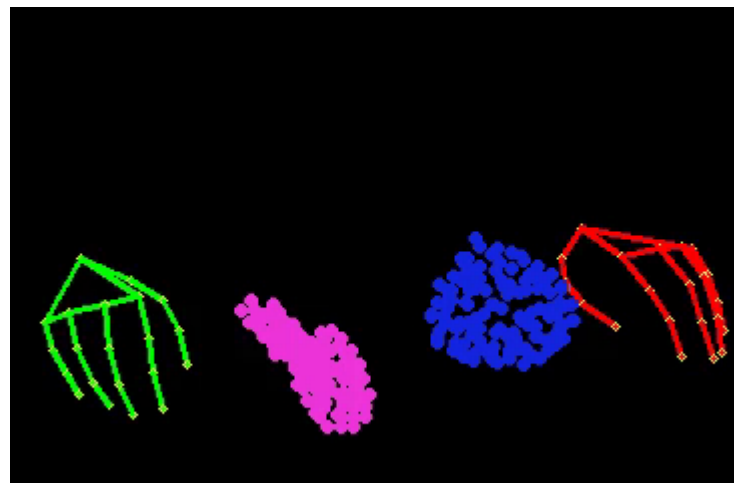
use spoon to put in bowl

Generated  
Motion

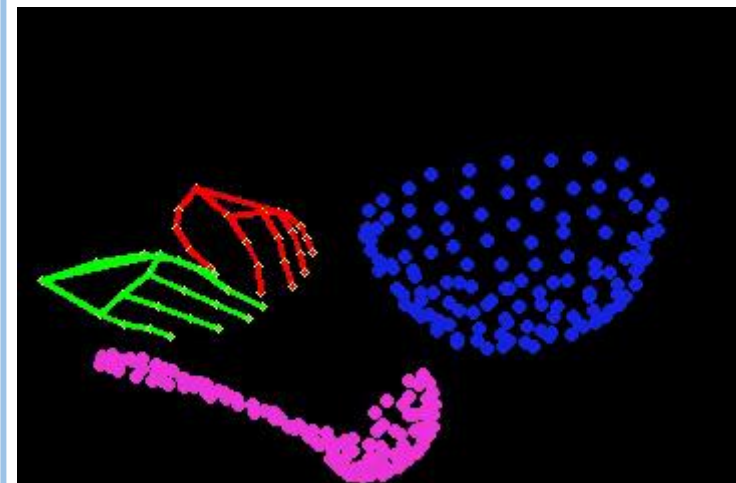
MDM



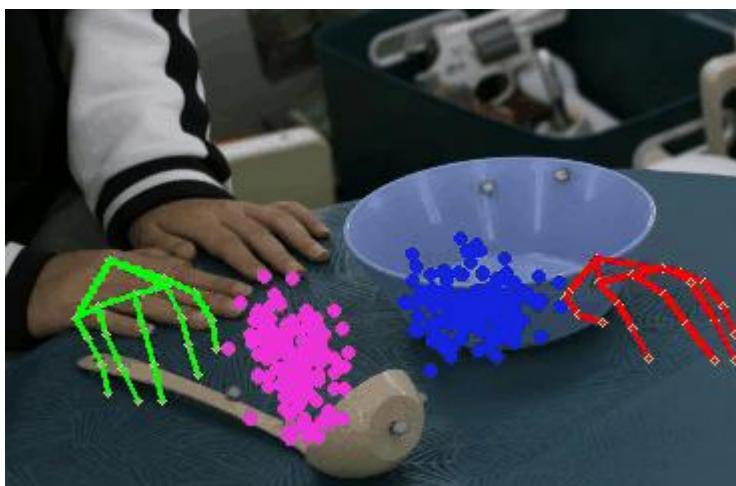
EMDM



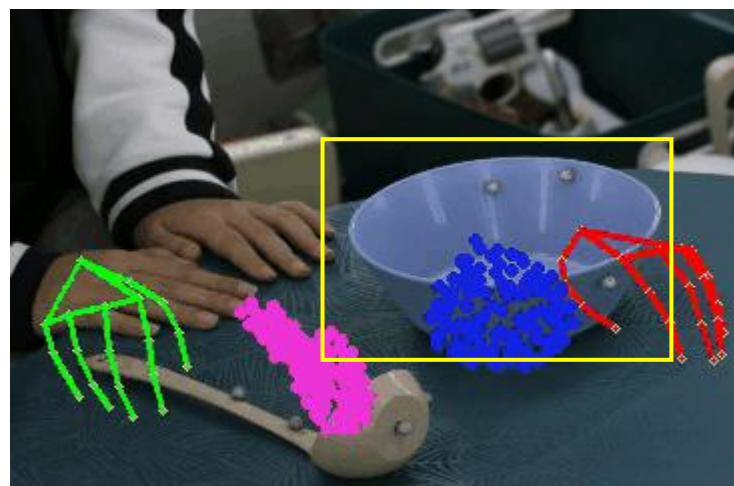
Ours



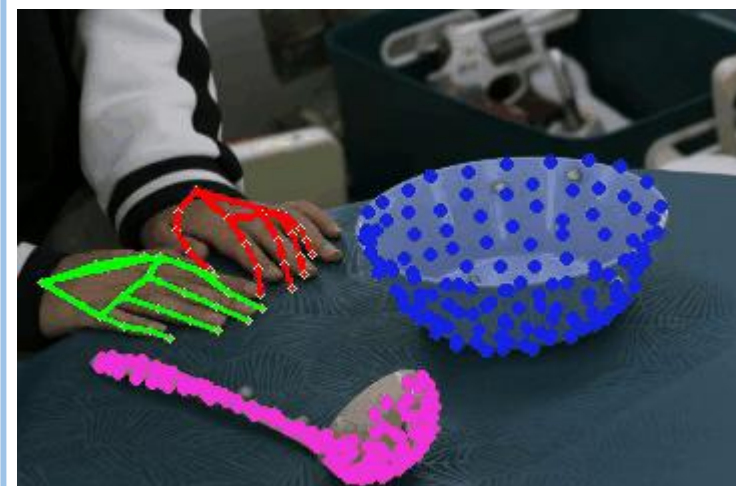
Overlaid  
Motion  
onto the Image



**Blurred outline, Low dynamic**



**Image-motion inconsistency**





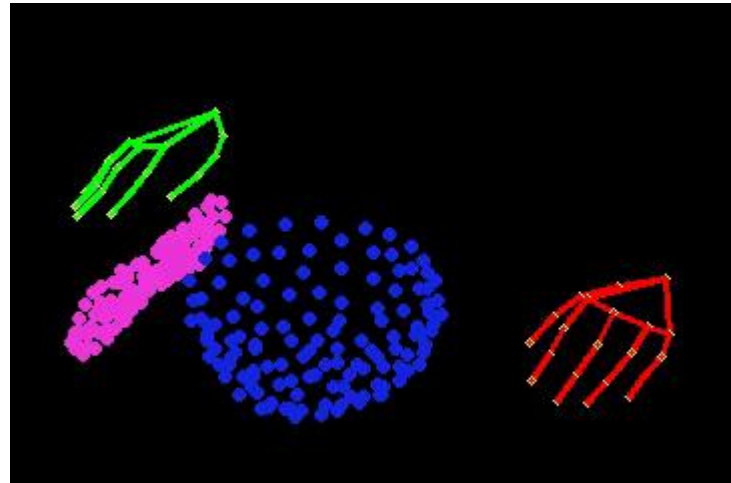
# 6. Generalization on the Real-World Data

use the spoon to scrape the bowl

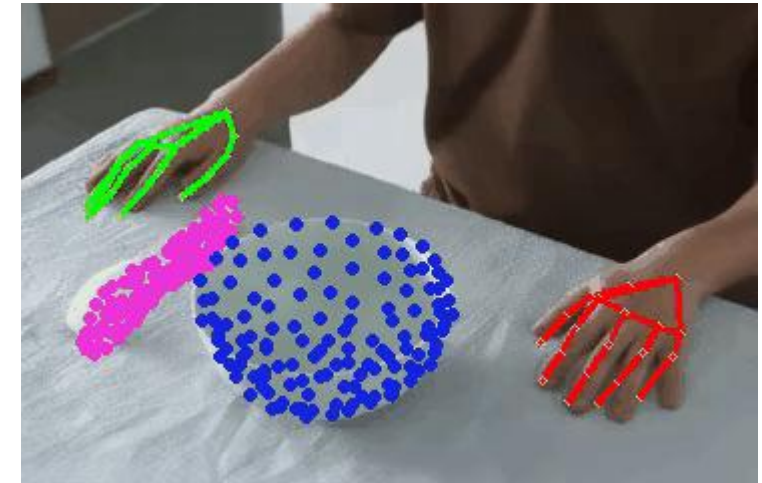
Generated Video



Generated Motion



Overlaid Results





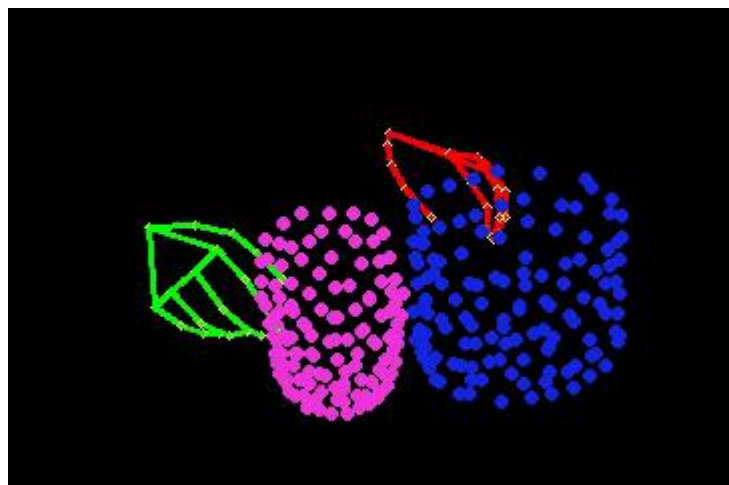
## 6. Generalization on the OAKINK2 Data

use the cup to pour into the bowl

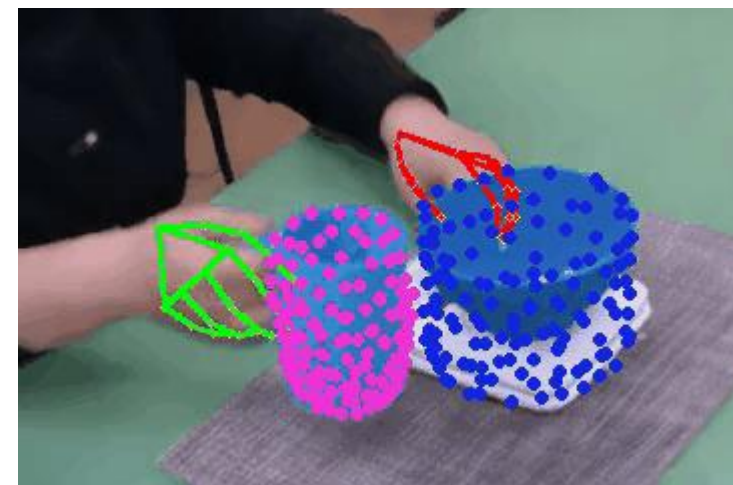
Generated Video



Generated Motion



Overlaid Results



Xinyu, Zhan, et al. "Oakink2: A dataset of bimanual hands-object manipulation in complex task completion."  
*Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2024.



# 7. QR Code for Our Project



[https://droliven.github.io/SViMo\\_project/](https://droliven.github.io/SViMo_project/)

**Thank you!**

## Acknowledgement

This work was supported by the National Natural Science Foundation of China (NSFC) No.62125107 and No.62272172.