

# CG-HOI: Contact-Guided 3D Human-Object Interaction Generation

CVPR 2024

Presenter: Jiahang Zhang  
2024.3.31

# Background: SMPL

---

## SMPL: A Skinned Multi-Person Linear Model

Matthew Loper<sup>\*1,2</sup> Naureen Mahmood<sup>†1</sup> Javier Romero<sup>†1</sup> Gerard Pons-Moll<sup>†1</sup> Michael J. Black<sup>†1</sup>

<sup>1</sup>Max Planck Institute for Intelligent Systems, Tübingen, Germany

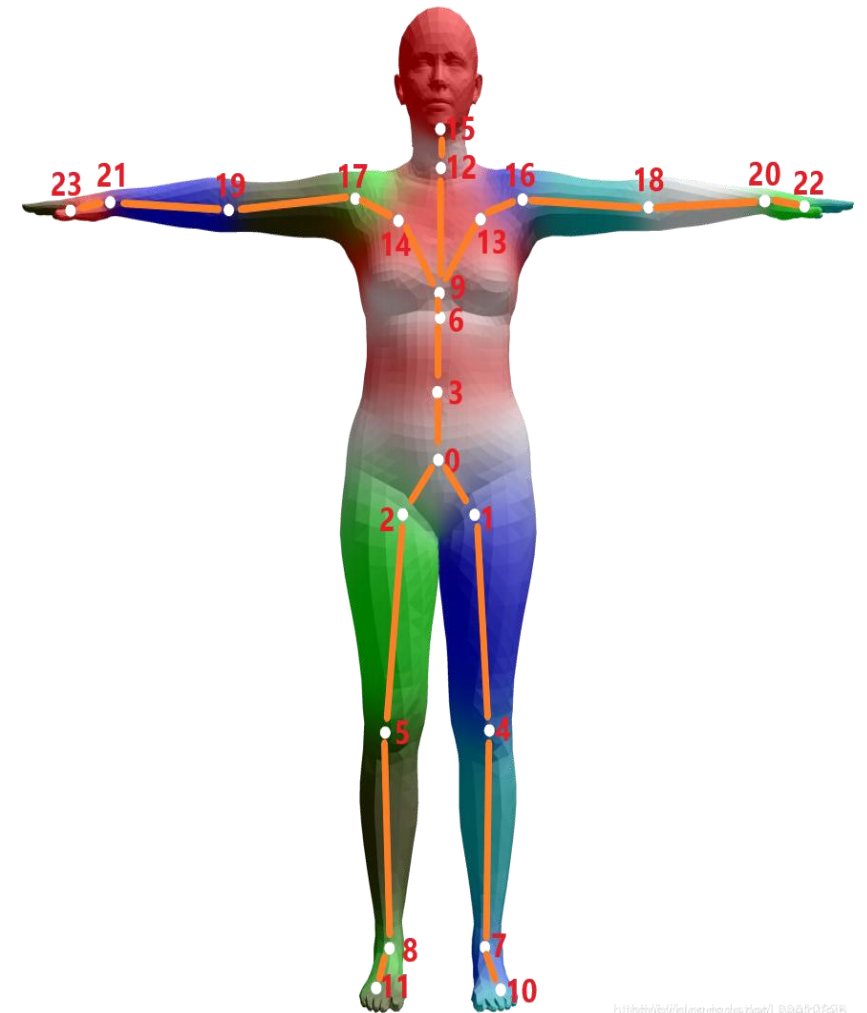
<sup>2</sup>Industrial Light and Magic, San Francisco, CA

### SMPL Model

- Shape Parameters (10,)
- Pose Parameters (24,3)
- Vertex Number  $N = 6890$

**SMPL-H:** hand+body

**SMPL-X:** face+hand+body



# Background: 3D Human Reconstruction

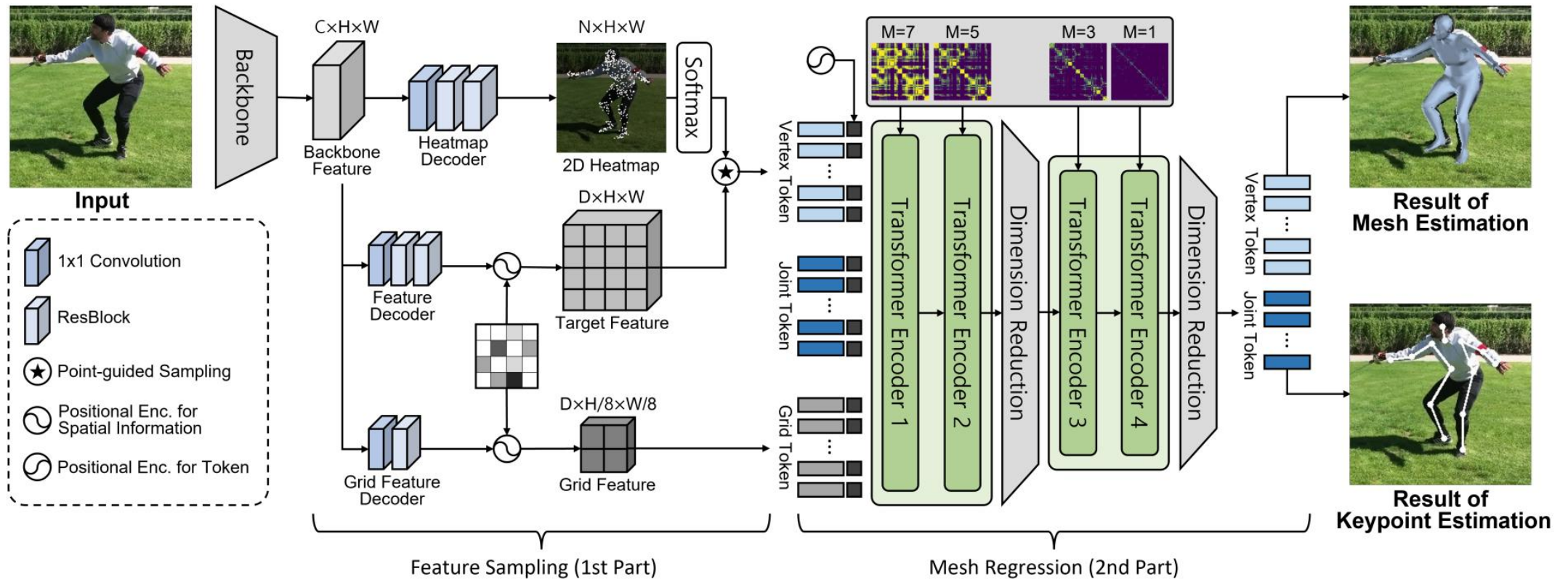
## Sampling is Matter: Point-guided 3D Human Mesh Reconstruction

Jeonghwan Kim<sup>1\*</sup> Mi-Gyeong Gwon<sup>1\*</sup> Hyunwoo Park<sup>1</sup>

Hyukmin Kwon<sup>2</sup> Gi-Mun Um<sup>2</sup> Wonjun Kim<sup>1†</sup>

<sup>1</sup>Konkuk University <sup>2</sup>Electronics and Telecommunications Research Institute



## Model-Free Methods



# Background: 3D Human Reconstruction

---

## PyMAF: 3D Human Pose and Shape Regression with Pyramidal Mesh Alignment Feedback Loop

Hongwen Zhang<sup>§†\*</sup>, Yating Tian<sup>†\*</sup>, Xinchu Zhou<sup>‡</sup>, Wanli Ouyang<sup>‡</sup>, Yebin Liu<sup>‡</sup>, Limin Wang<sup>†</sup> , Zhenan Sun<sup>§</sup> 

<sup>§</sup>CRIPAC, NLPR, Institute of Automation, Chinese Academy of Sciences, China

<sup>†</sup>State Key Laboratory for Novel Software Technology, Nanjing University, China

<sup>‡</sup>The University of Sydney, Australia <sup>‡</sup>Department of Automation, Tsinghua University, China

## Model-Based Methods

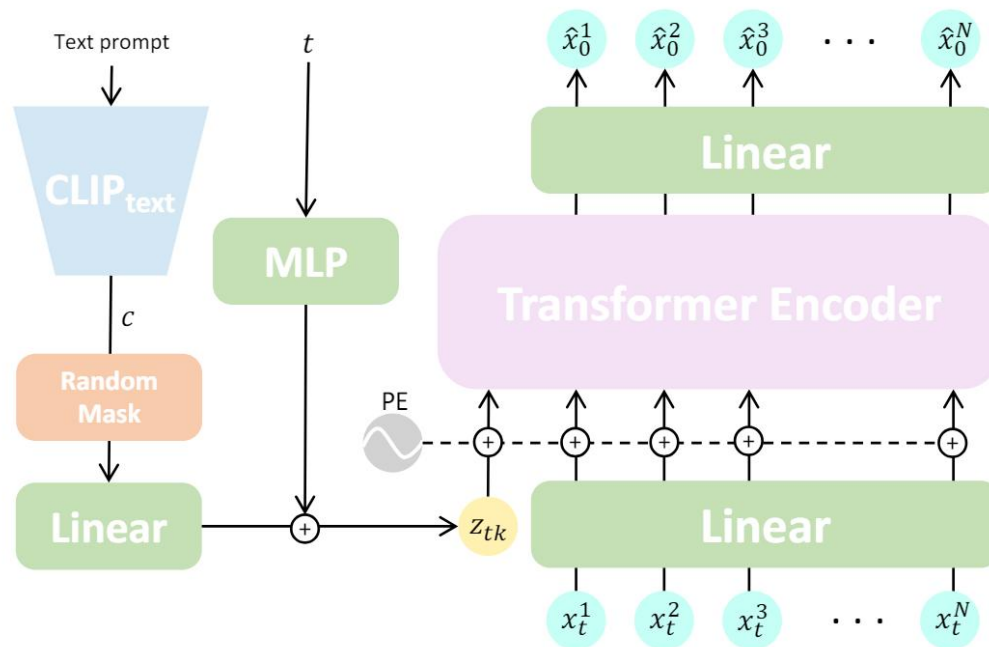
- Pose, shape, and camera parameters  $\Theta$

$$\mathcal{L}_{reg} = \lambda_{2d} \|K - \hat{K}\|^2 + \lambda_{3d} \|J - \hat{J}\|^2 + \lambda_{para} \|\Theta - \hat{\Theta}\|^2, \quad (2)$$

where  $\|\cdot\|^2$  is the squared L2 norm,  $\hat{K}$ ,  $\hat{J}$ , and  $\hat{\Theta}$  denote the ground truth 2D keypoints, 3D joints, and model parameters, respectively.

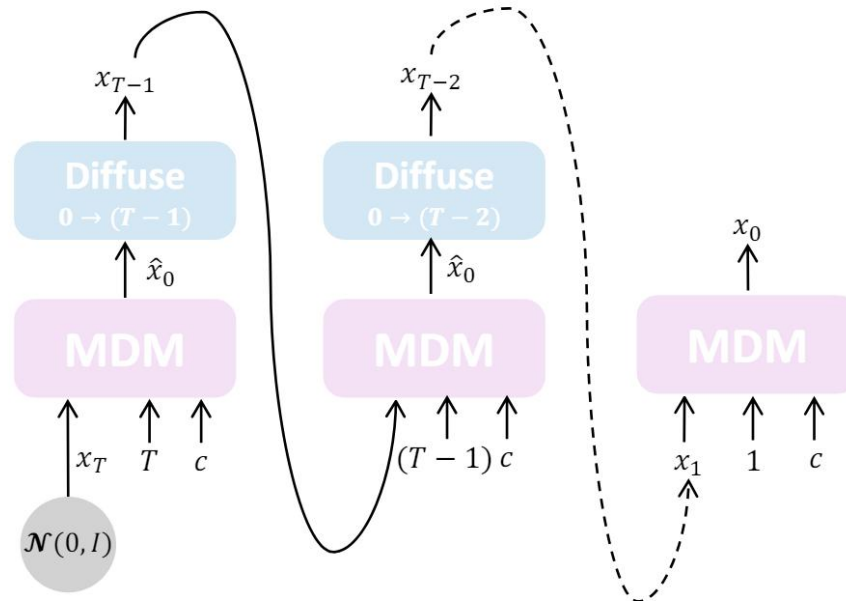
# Background: 3D Human Motion Generation

## Diffusion-Based



## HUMAN MOTION DIFFUSION MODEL

Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafir,  
Daniel Cohen-Or and Amit H. Bermano  
Tel Aviv University, Israel  
guytevet@mail.tau.ac.il

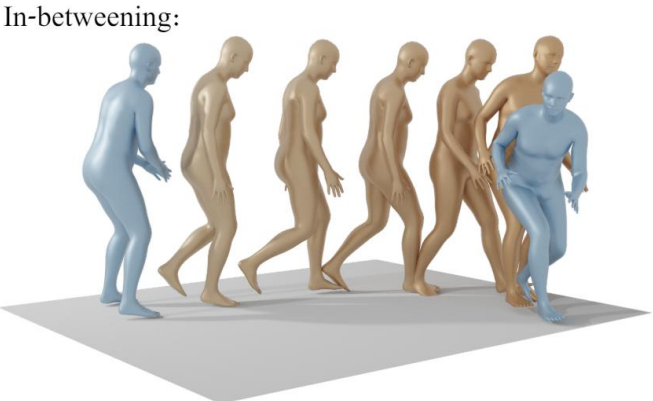


# Background: 3D Human Motion Generation

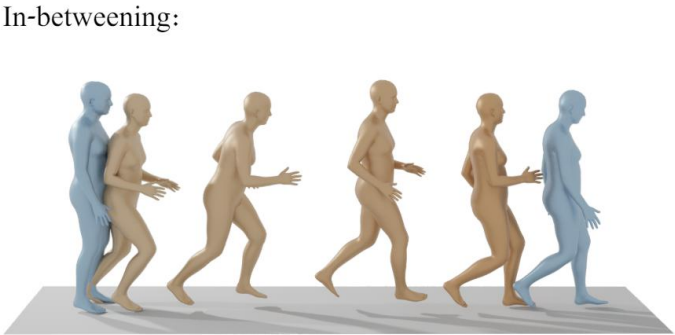
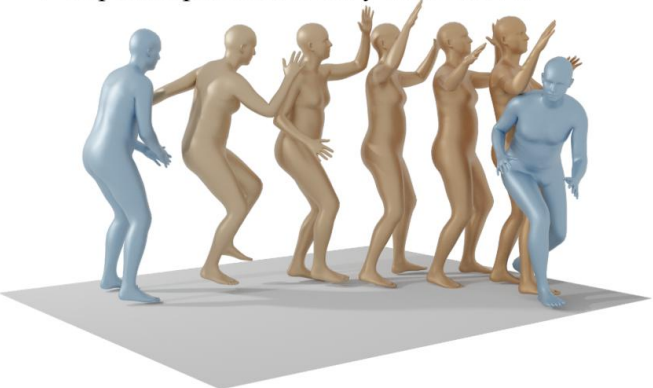
## HUMAN MOTION DIFFUSION MODEL

Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafir,  
Daniel Cohen-Or and Amit H. Bermano  
Tel Aviv University, Israel  
guytevet@mail.tau.ac.il

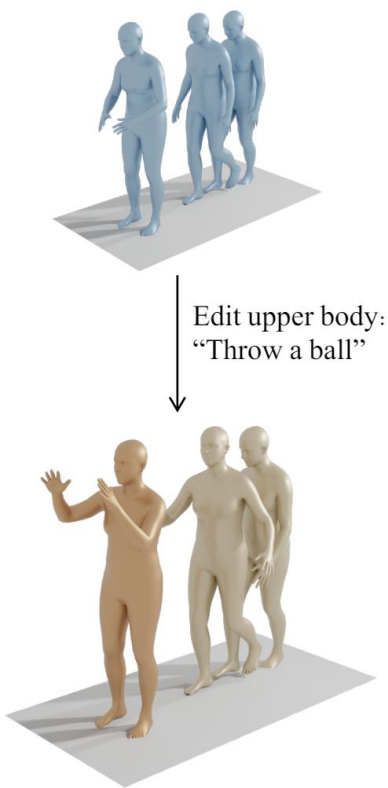
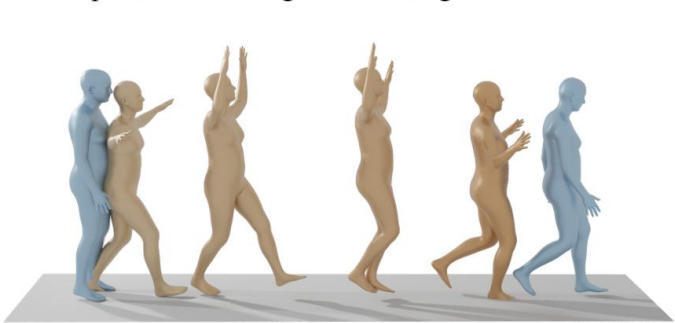
### Diffusion-Based



+ "A person performs a crazy dance move."



+ "A person is walking while raising hands."





# Background: 3D Human Motion Generation

## Guided Motion Diffusion for Controllable Human Motion Synthesis

Korraue Karunratanakul<sup>1</sup> Konpat Preechaku<sup>2</sup> Supasorn Suwajanakorn<sup>2</sup> Siyu Tang<sup>1</sup>

<sup>1</sup>ETH Zürich, Switzerland <sup>2</sup>VISTEC, Thailand

<https://korraue.github.io/gmd-project/>

## Diffusion-Based



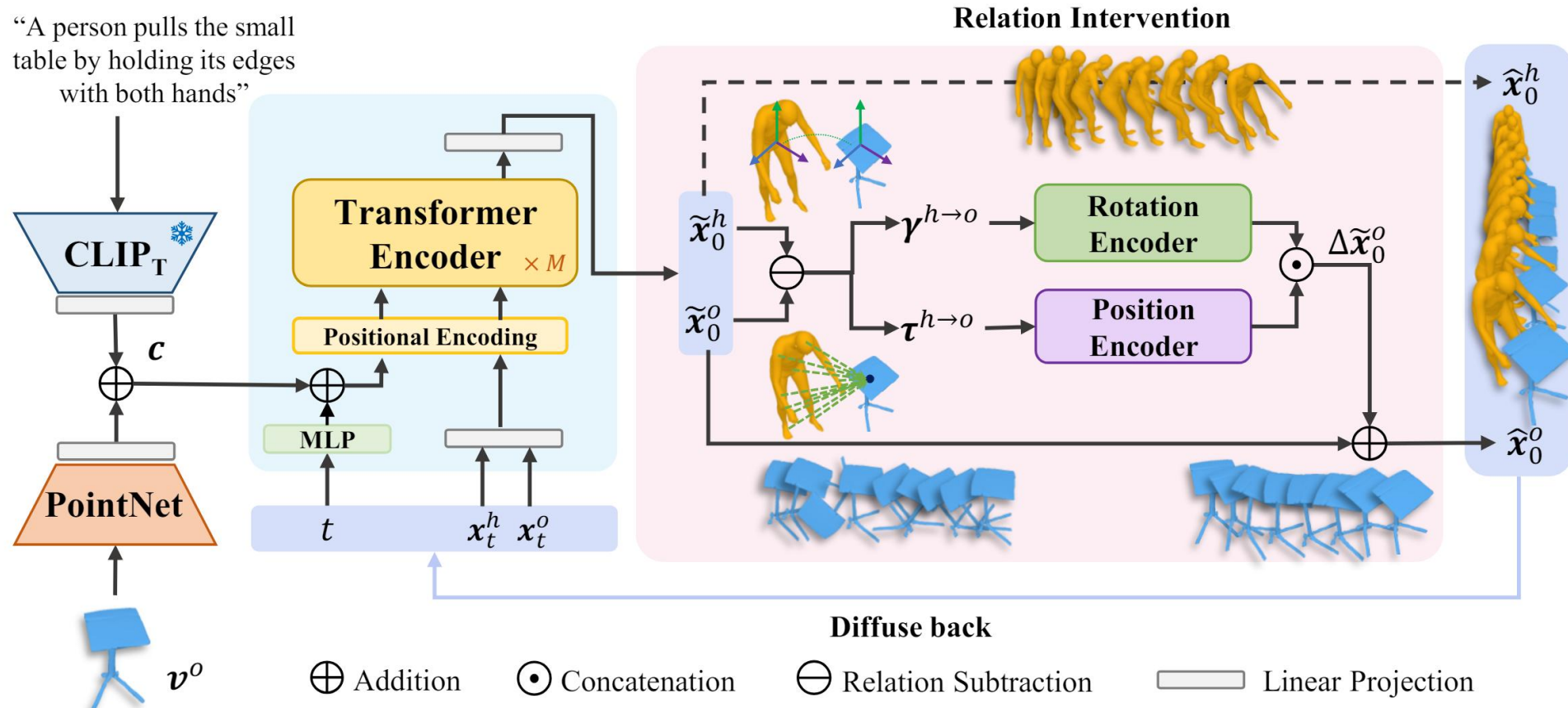
# Background: HOI Generation

## THOR: Text to Human-Object Interaction Diffusion via Relation Intervention

Qianyang Wu<sup>1</sup>, Ye Shi<sup>1</sup>, Xiaoshui Huang<sup>2</sup>, Jingyi Yu<sup>1</sup>, Lan Xu<sup>1</sup>, and Jingya Wang<sup>1</sup>

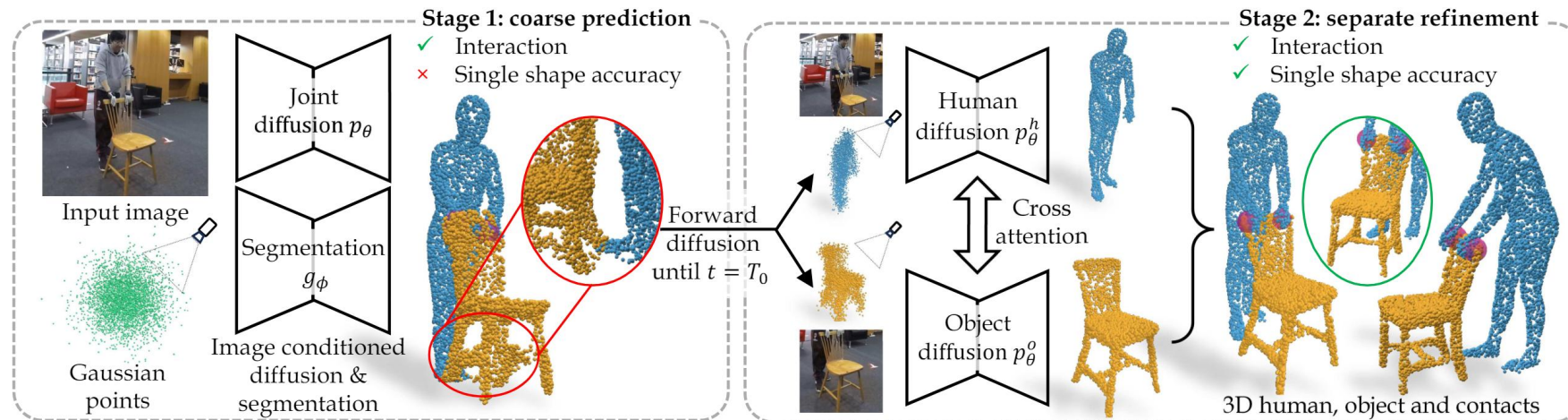
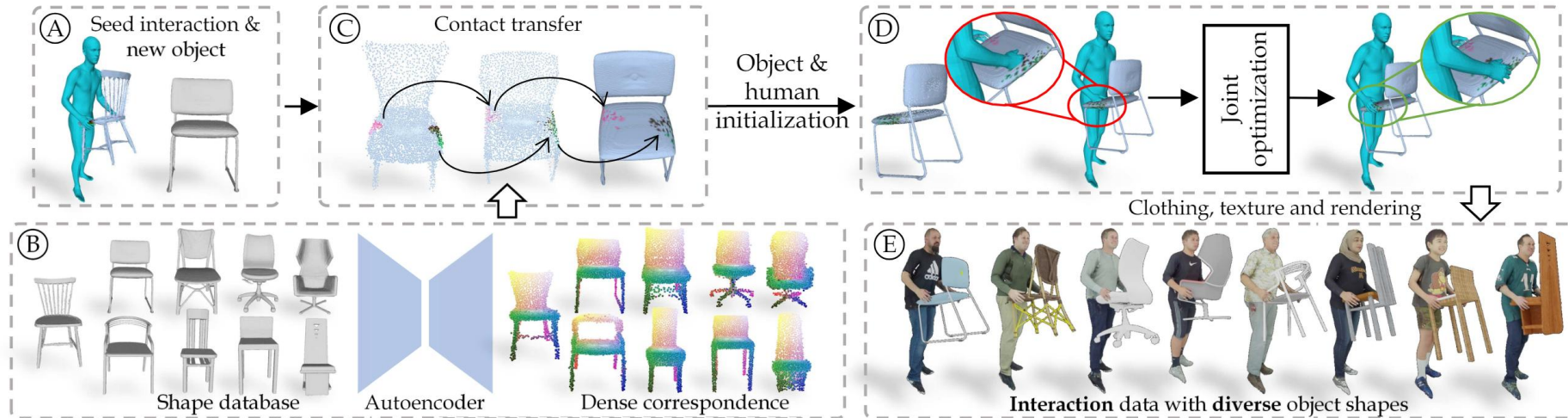
<sup>1</sup> ShanghaiTech University

<sup>2</sup> Shanghai AI Laboratory





# Background: HOI Generation

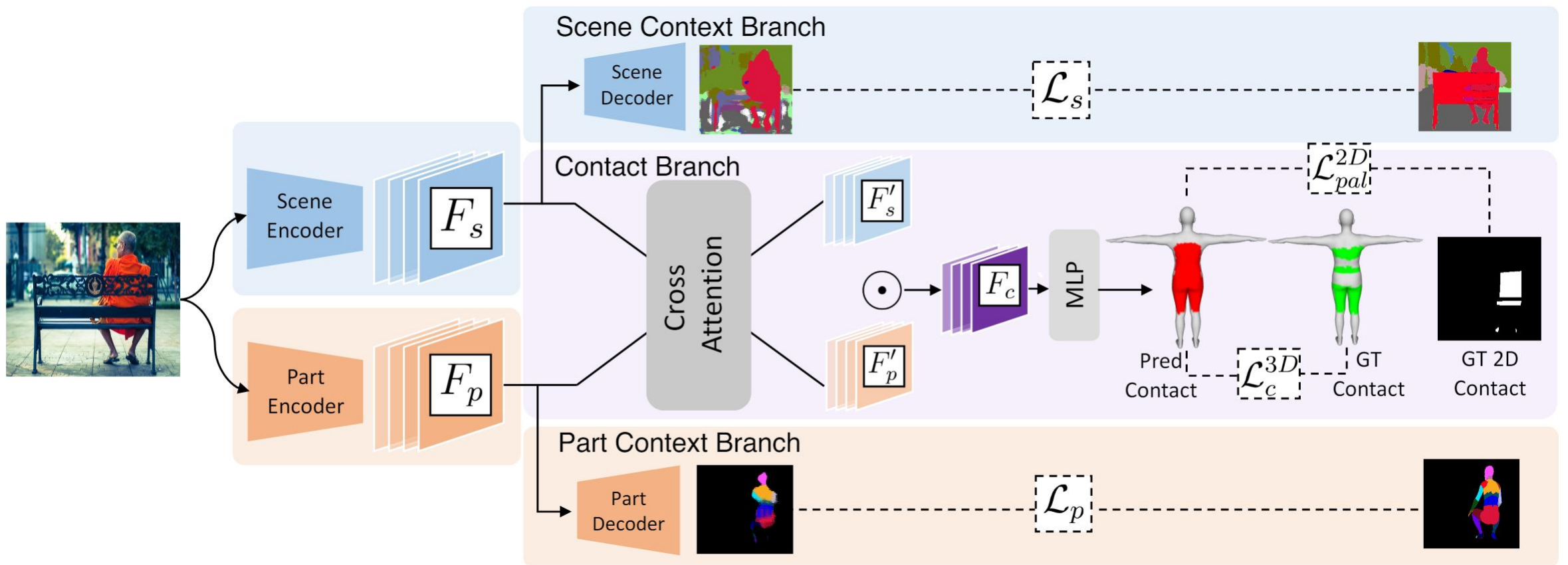


# Background: Contact Prediction for HOI

## DECO: Dense Estimation of 3D Human-Scene Contact In The Wild

Shashank Tripathi<sup>1\*</sup> Agniv Chatterjee<sup>1\*</sup> Jean-Claude Passy<sup>1</sup> Hongwei Yi<sup>1</sup>  
Dimitrios Tzionas<sup>2</sup> Michael J. Black<sup>1</sup>

<sup>1</sup>Max Planck Institute for Intelligent Systems, Tübingen, Germany <sup>2</sup>University of Amsterdam, the Netherlands  
{stripathi, a chatterjee, jpassy, hyi, black}@tue.mpg.de d.tzionas@uva.nl

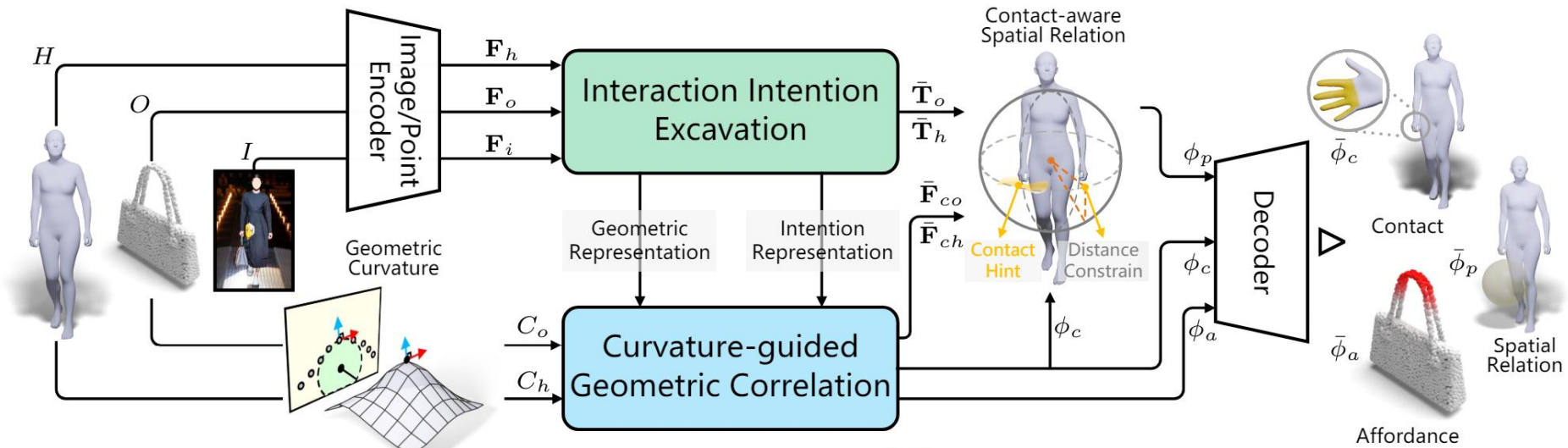
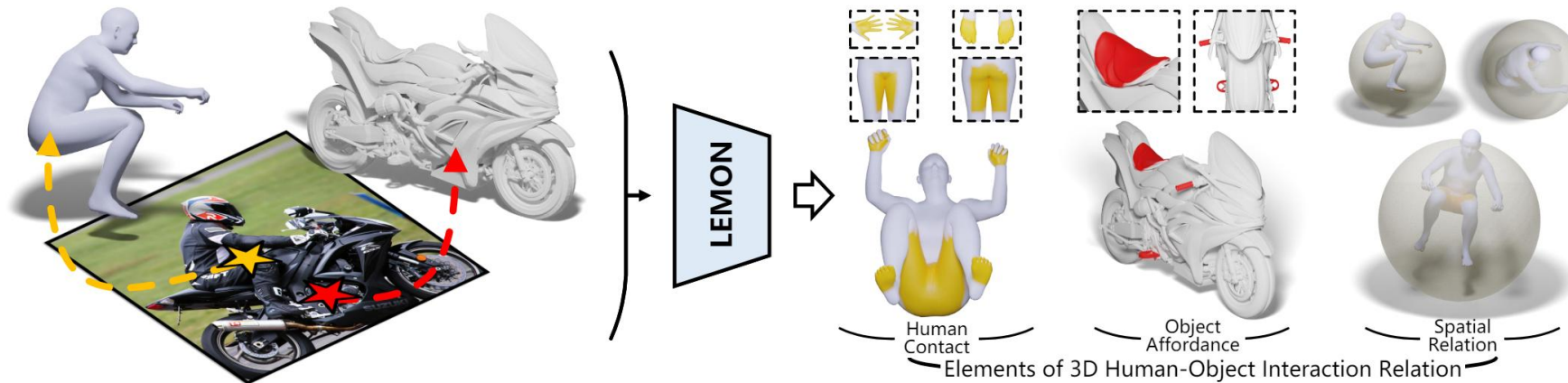


# Background: Contact Prediction for HOI

Yuhang Yang<sup>1</sup>, Wei Zhai<sup>1,†</sup>, Hongchen Luo<sup>1</sup>, Yang Cao<sup>1,2</sup>, Zheng-Jun Zha<sup>1</sup>

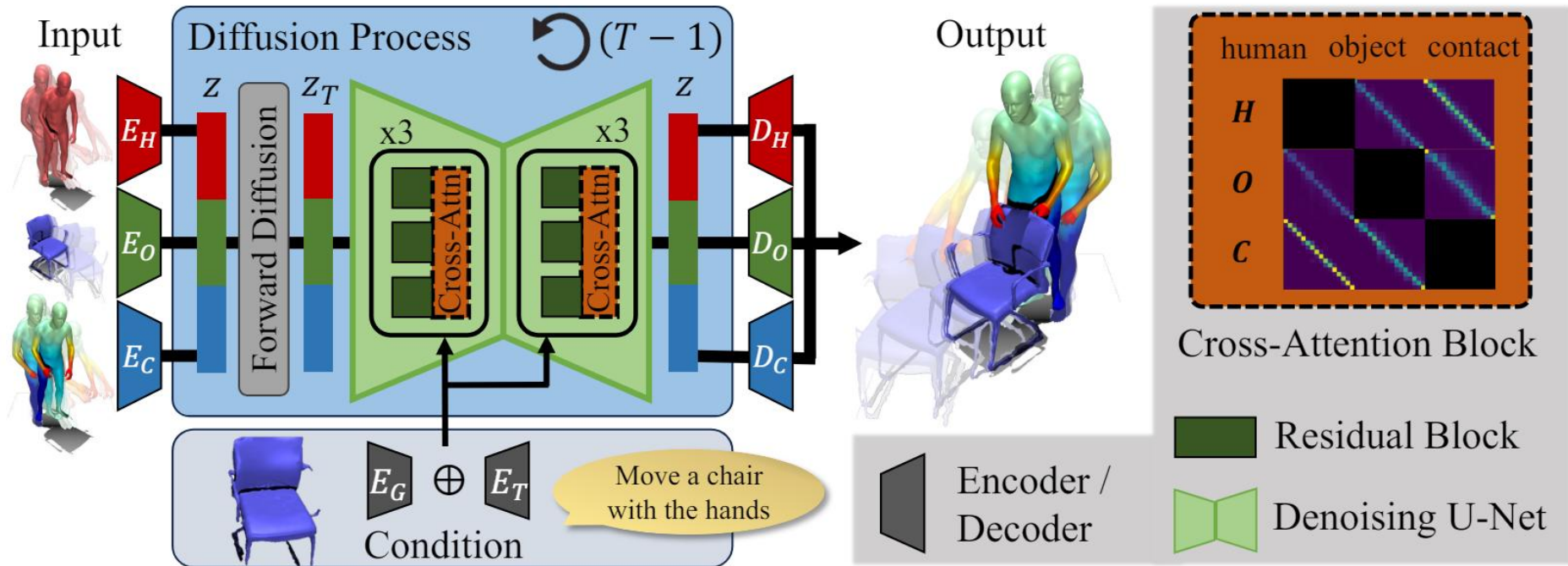
<sup>1</sup> University of Science and Technology of China

<sup>2</sup> Institute of Artificial Intelligence, Hefei Comprehensive National Science Center





# CG-HOI: Method

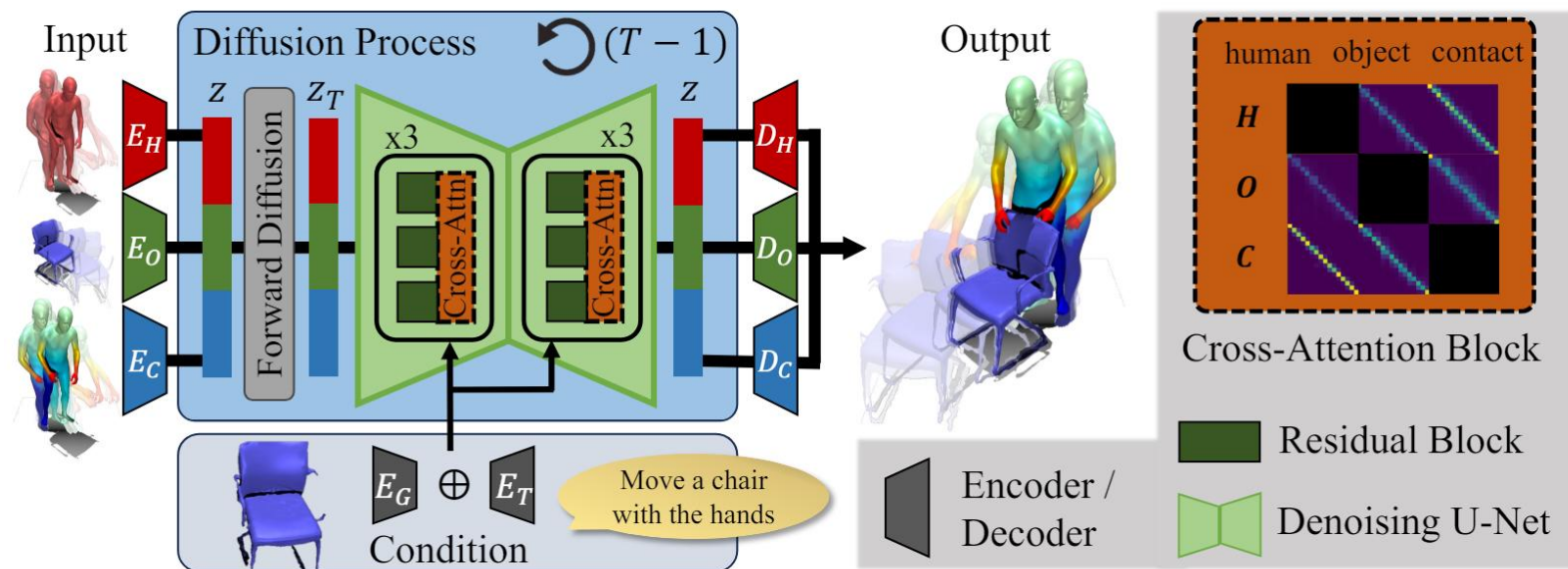


- Diffusion-based
- Conditioned on text description  $T$ , static geometry  $G$
- Contact-guided inference

# CG-HOI: Method

## Diffusion Process

- Human surface: SMPL
- Object Transformation: global translation, rotation matrix
- Contact: distance between markers to the closest object geometry point

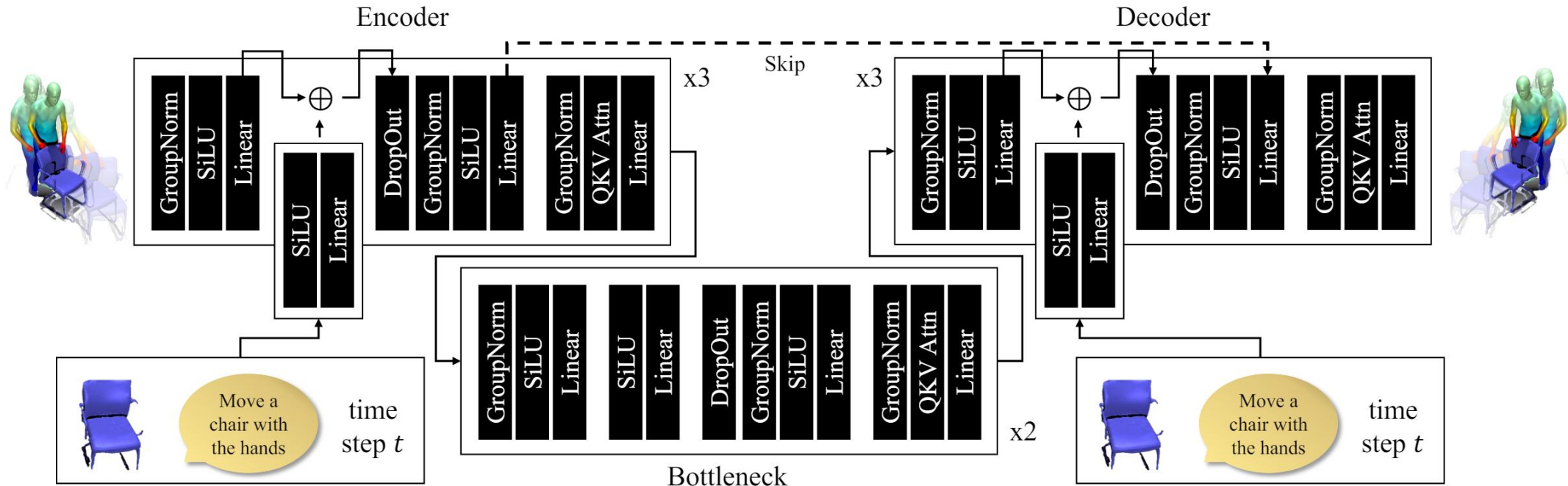




# CG-HOI: Method

## Diffusion Process

- Text encoder: CLIP
- Object Geometry encoder: PointNet
- Cross-Attention Mechanism for  $h_i, o_i, c_i$ , to model inter-dependency



### Object Motion Modeling

- *Object motion is naturally most influenced by parts of the human body in very close contact to the object.*
- Predict one hypotheses for each marker, instead one for the whole seq.

Formally, we predict object transformation hypotheses  $o_i^j$  for each contact point on the human body, and weigh them with the inverse of their predicted contact distance  $c_i^j$ :

$$o_i = \frac{1}{\sum_j c_i} \sum_{j=0}^N (\max(|c_i|) - |c_i^j|) o_i^j, \quad (5)$$



### Loss Formulation

- Between Prediction and Ground-Truth
- Classifier-Free Guidance

$$\mathbf{L} = \lambda_h ||h_i - \hat{h}_i||_1 + \lambda_o ||o_i - \hat{o}_i||_1 + \lambda_c ||c_i - \hat{c}_i||_2,$$

### Interaction Generation

- Explicit constraints during inference on human-object contact
- Apply a a cost function  $G((x)_t)$  at each time step

$$\hat{\mu}_t = \mu_t + s \sum_t \nabla_{x_t} G(x_t),$$

- The input can also be replaced by the object trajectory without any new training.

# CG-HOI: Experiment

---

## Dataset

- CHAIRS: 46 subjects as their SMPL-X, interaction with chairs and sofas
- BEHAVE: 8 participants as their SMPL-H alongside 20 different objects

## Evaluation Metrics

- **R-Precision**: the closeness of the text condition and generated HOI in latent feature space.
- **FID**: the similarity between generated and ground-truth distribution in encoded feature space.
- **Diversity and MultiModality**: motion variance across all text descriptions.
- **Perceptual User Study**



# CG-HOI: Experiment

## Quantitative Result

		BEHAVE			CHAIRS	
Task	Approach	R-Prec. (top-3) $\uparrow$	FID $\downarrow$	Diversity $\rightarrow$		
	Real (human)	0.73	0.09	4.23		
Text-Cond. Human Only	MDM [71]	0.52	4.54	5.44		
	InterDiff [84]	0.49	5.36	3.98		
	<b>Ours</b>	<b>0.60</b>	<b>4.26</b>	<b>4.92</b>		
	Real	0.81	0.17	6.80		
Motion-Cond. HOI	InterDiff [84]	0.68	3.86	5.62		
	<b>Ours</b>	<b>0.71</b>	<b>3.52</b>	<b>6.89</b>		
Text-Cond. HOI	MDM [71]	0.49	9.21	6.51		
	InterDiff [84]	0.53	8.70	3.85		
	<b>Ours</b>	<b>0.62</b>	<b>6.31</b>	<b>6.63</b>		

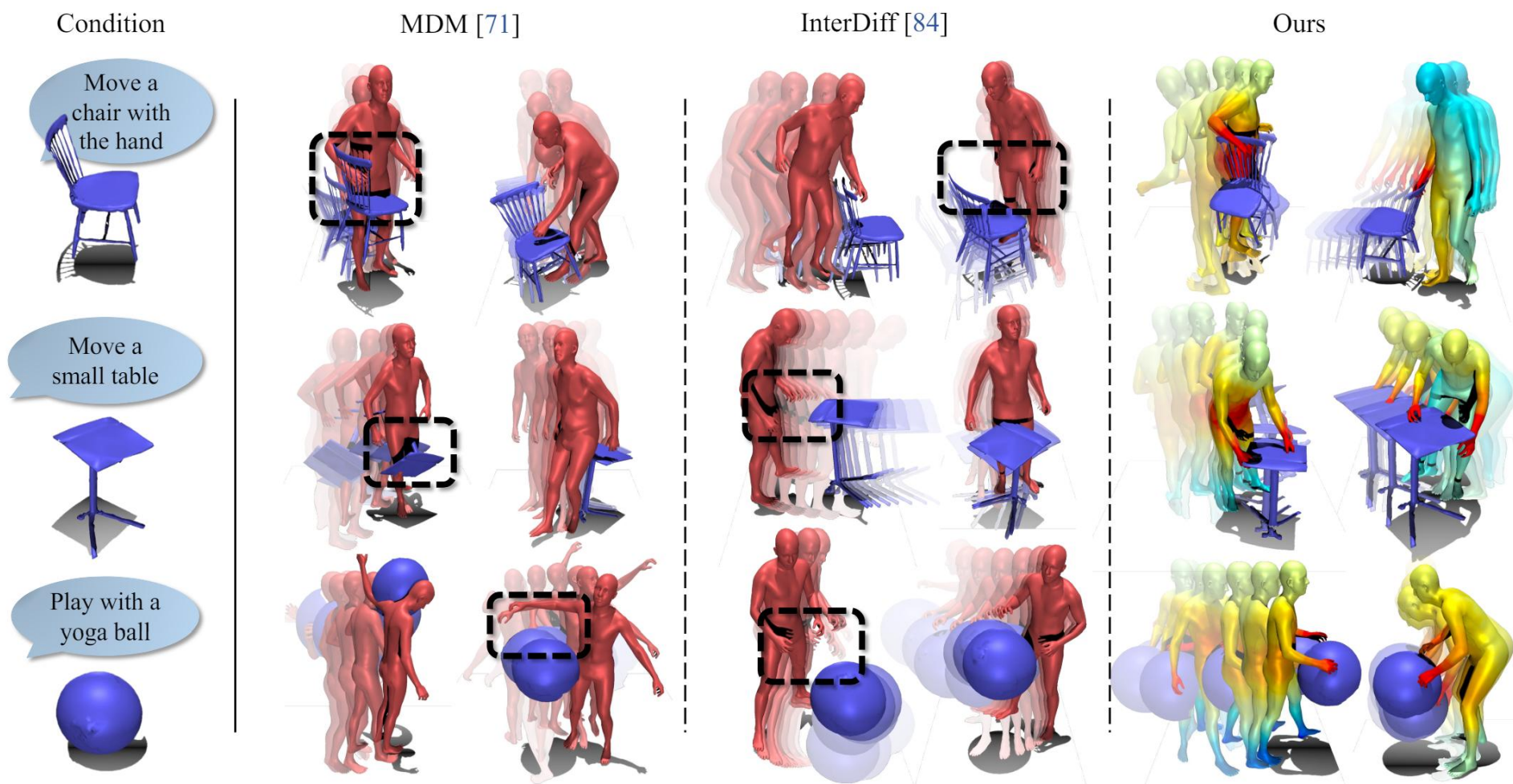
Method	Realistic Interactions (%)	Text Coherence (%)
MDM	81.8%	79.5%
InterDiff	72.8%	73.1%
GT	34.3%	48.7%

**Table 1.** Quantitative comparison with state-of-the-art approaches, on the human pose sequence, and motion-cond. denotes predictions of object behavior. For metrics with  $\rightarrow$ , results closer to the real distribution are better. Our approach outperforms these baselines in all three settings, indicating a strong learned correlation between human and object motion.

**Figure 5.** Perceptual User Study. Participants significantly favor our method over baselines, for overall realism and text coherence.

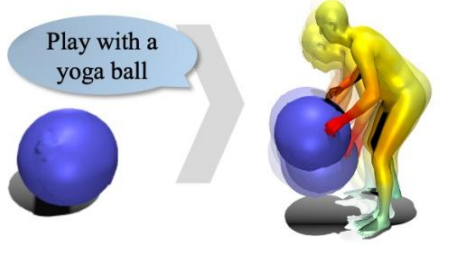
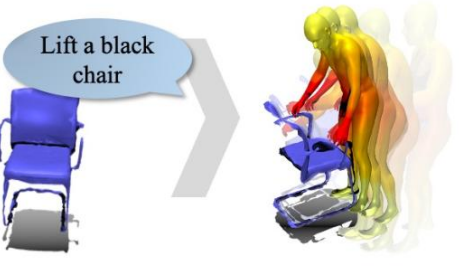
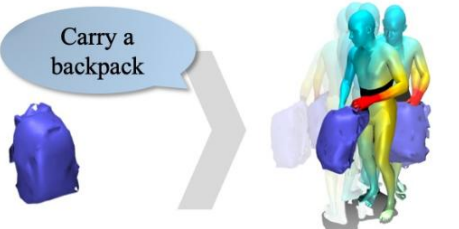
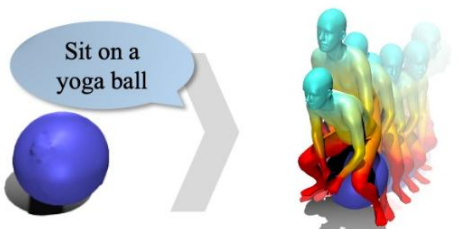
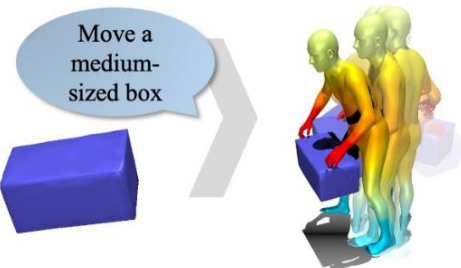
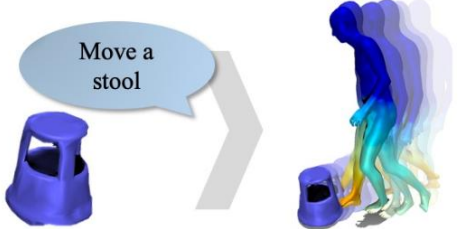
# CG-HOI: Experiment

## Qualitative Result



# CG-HOI: Experiment

## Qualitative Result

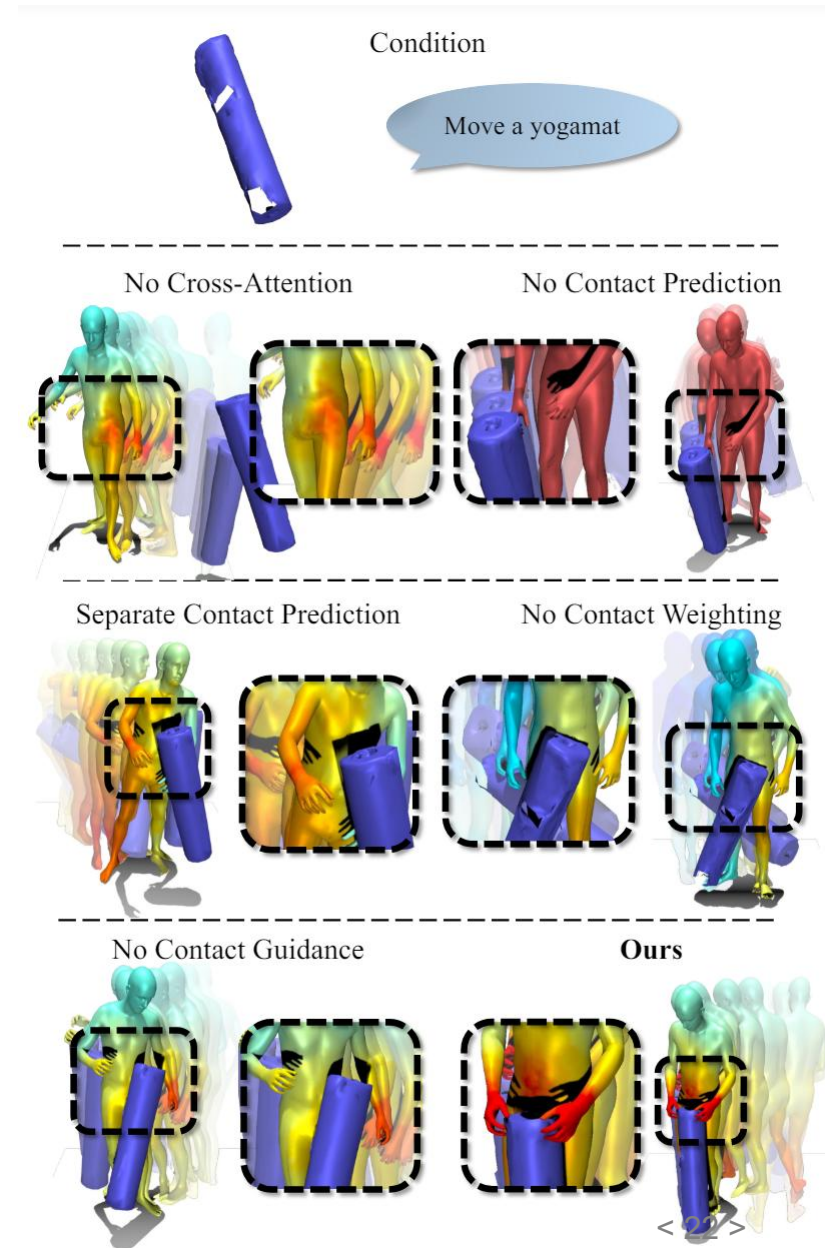




# CG-HOI: Experiment

## Ablation Study

Approach	BEHAVE				CHAIRS			
	R-Prec. (top-3) ↑	FID ↓	Diversity →	MModality →	R-Prec. (top-3) ↑	FID ↓	Diversity →	MModality →
Real	0.81	0.17	6.80	6.24	0.87	0.02	9.91	6.12
No cross-attention	0.35	10.44	8.23	7.40	0.49	10.84	12.22	10.64
No contact prediction	0.41	9.64	10.10	6.89	0.41	8.53	11.56	9.15
Separate contact pred.	0.47	8.01	5.12	5.12	0.52	9.34	7.65	4.62
No contact weighting	0.55	8.54	6.52	5.29	0.64	7.55	8.56	5.45
No contact guidance	0.59	7.22	7.84	5.30	0.70	7.41	8.05	5.76
<b>Ours</b>	<b>0.62</b>	<b>6.31</b>	<b>6.63</b>	<b>5.47</b>	<b>0.74</b>	<b>6.74</b>	<b>8.91</b>	<b>5.94</b>

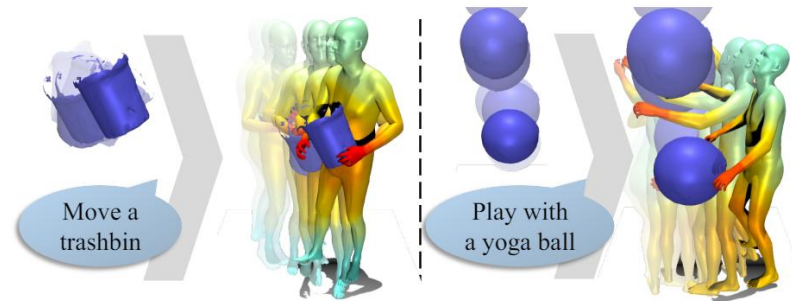


# CG-HOI: Experiment

---

## Some Applications

- Human motion generation given object trajectory.



**Figure 7.** Given an object trajectory at inference time, our method can generate corresponding human motion without re-training.

- Populating 3D scans. Generate realistic human motion sequences given a static scene.





# Conclusion

---

- Contact guided method for diffusion-based HOI generation.
- The first approach to address the task of generating realistic 3D HOI from text.
- Can generate motion that lasts up to 3 seconds.
- Seem to lack the temporal consistency constraint?
- 3D HOI training data is costly.

# STRUCT Group Seminar

---

Thanks!