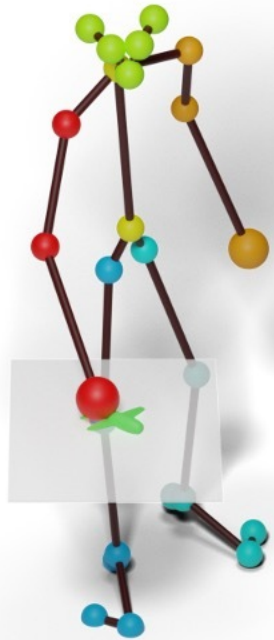# 3D Human Behavior Generation through Action and Interaction Synthesis
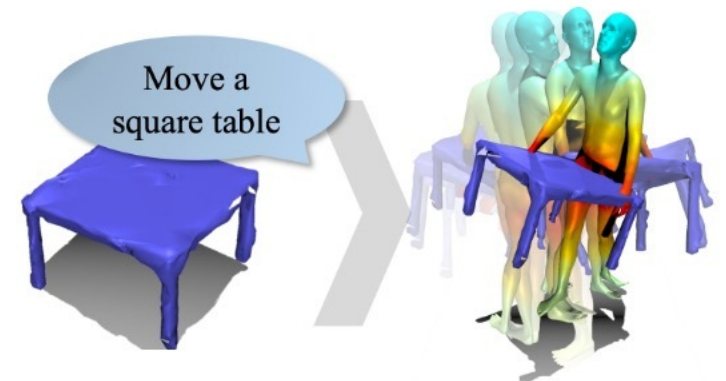
**Christian Diller**

**Supervisor: Prof. Angela Dai**
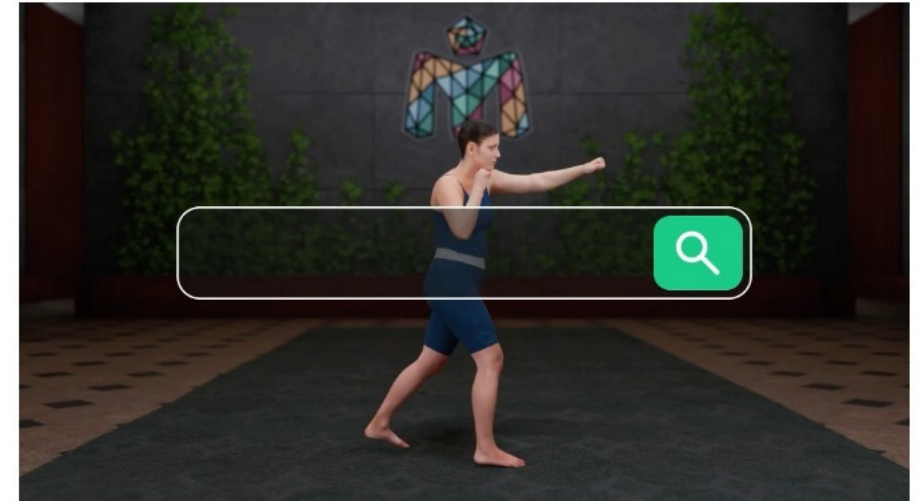
**Tuesday, 10th December 2024**

# Motivation: Understanding Human Behavior in 3D

- **Human behavior understanding is important for perception**

  - **Higher-order understanding of human-machine interaction**

  - **Anticipatory action vs. perceptual reaction**



- **Human environments are made by humans for humans**

- **Human motion generation in 3D**

  - **Allows for more fine-grained actions, e.g., grasping objects**

  - **Enables direct interactions with an environment**
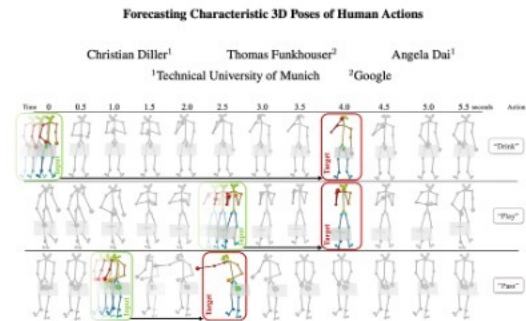


Move a
square table

# Applications



- **Human-centered assistive systems**
  - **Interaction between humans and robots in a shared physical space**
  - **Assistance robotics in medicine and care**



- **Autonomous driving**
  - **Forecasting interactions between cars & pedestrians**

- **Content Creation**
  - **Plausible human motion from sparse input (e.g., text)**

# 3D Human Behavior Generation: Action & Interaction

## Efficient Action Representation



**Forecasting Characteristic 3D Poses [1]**

## Complex Action Sequences



**FutureHuman3D [2]**

## Human-Object Interactions



**CG-HOI [3]**

[1] Diller, Christian, Thomas Funkhouser, and Angela Dai. "Forecasting characteristic 3d poses of human actions." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022.

[2] Diller, Christian, Thomas Funkhouser, and Angela Dai. "FutureHuman3D: Forecasting Complex Long-Term 3D Human Behavior from Video Observations." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024.

[3] Diller, Christian, and Angela Dai. "Cg-hoi: Contact-guided 3d human-object interaction generation." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024.

# 3D Human Behavior Generation: Action & Interaction

## Efficient Action Representation



**Forecasting Characteristic 3D Poses [1]**

## Complex Action Sequences



**FutureHuman3D [2]**

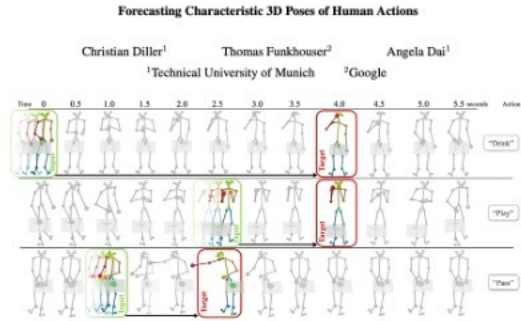## Human-Object Interactions



**CG-HOI [3]**

[1] Diller, Christian, Thomas Funkhouser, and Angela Dai. "Forecasting characteristic 3d poses of human actions." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022.

[2] Diller, Christian, Thomas Funkhouser, and Angela Dai. "FutureHuman3D: Forecasting Complex Long-Term 3D Human Behavior from Video Observations." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024.

[3] Diller, Christian, and Angela Dai. "Cg-hoi: Contact-guided 3d human-object interaction generation." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024.

# Forecasting Characteristic 3D Poses of Human Actions

How to efficiently represent 3D human motion sequences?

**Christian Diller**        **Thomas Funkhouser**        **Angela Dai**

TI.ITI        G        TI.ITI

# Time-Based Future Human Motion Prediction



$x_0$    $x_1$    $x_2$    $x_3$    $x_4$    $x_5$    $x_6$    $x_7$    $x_8$    $x_9$

Joint Location

Poses at fixed time steps

Continuous Movement

Time

# Forecasting Characteristic 3D Human Poses

# Task: Characteristic 3D Poses for Action Goals



Observation → Action Goal

Input → Target

Observation → Action Goal

Input → Target

# Dataset: Characteristic 3D Poses on GRAB [1]



Input

"Drink"

Target

**Original GRAB [1] Dataset** | **3D Skeleton Sequence** | **Pose Annotations**

[1] Taheri, O., Nima Ghorbani, Michael J. Black and Dimitrios Tzionas. "GRAB: A Dataset of Whole-Body Human Grasping of Objects." ECCV (2020).

# Dataset: Characteristic 3D Poses on Human3.6m [1]



"Phoning"

Characteristic Pose

[1] Ionescu, Catalin, et al. "Human3. 6m: Large scale datasets and predictive methods for 3d human sensing in natural environments." IEEE transactions on pattern analysis and machine intelligence 36.7 (2013): 1325-1339.

# Method: Architecture



Input Pose Sequence

Per-Voxel Offsets

k Skeleton Samples

Offsets

Encoder

Attention

Heatmap

Previous Joint Predictions (if any)

Multi-Modal Heatmap

Sampling

k Skeleton Samples (without offsets)

# Method: Autoregressive Prediction



Right Hand Prediction → Left Hand Prediction → Body Prediction → Pose Refinement

# Method: Pose Refinement

- **End-Effector Locations**

- **Bone-Lengths, as observed in input**

- **Joint angles, as observed in input**

- **Heatmap joint probability**

$$\mathbf{E}_R(\mathbf{x}, \mathbf{e}, \mathbf{b}, \theta, H) =$$



**Initial Prediction**            **Refined**

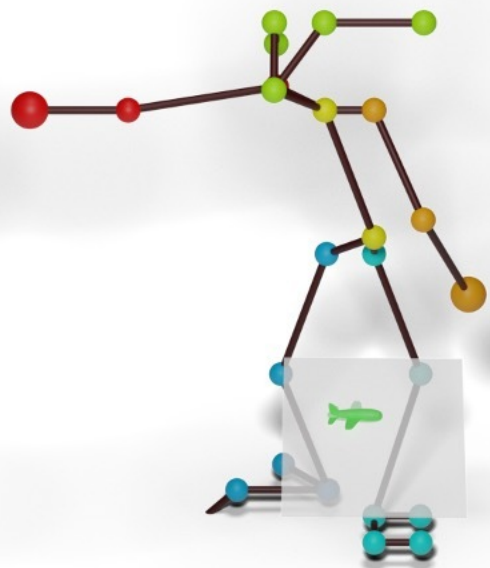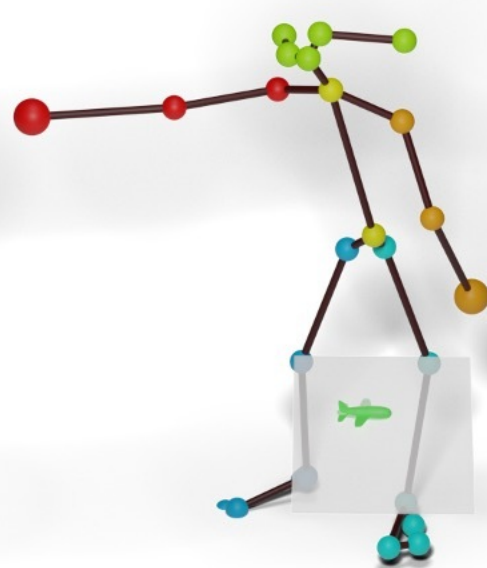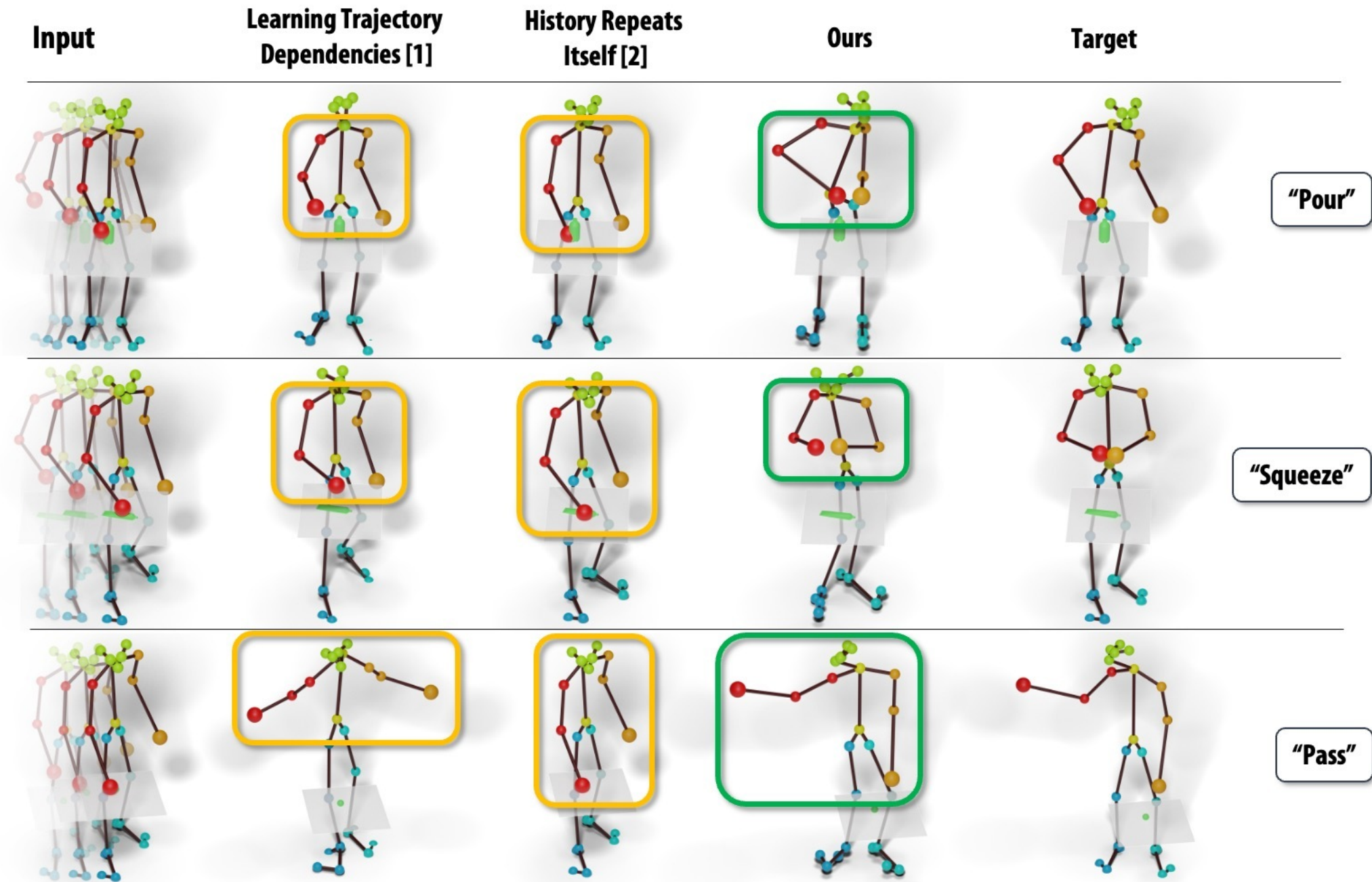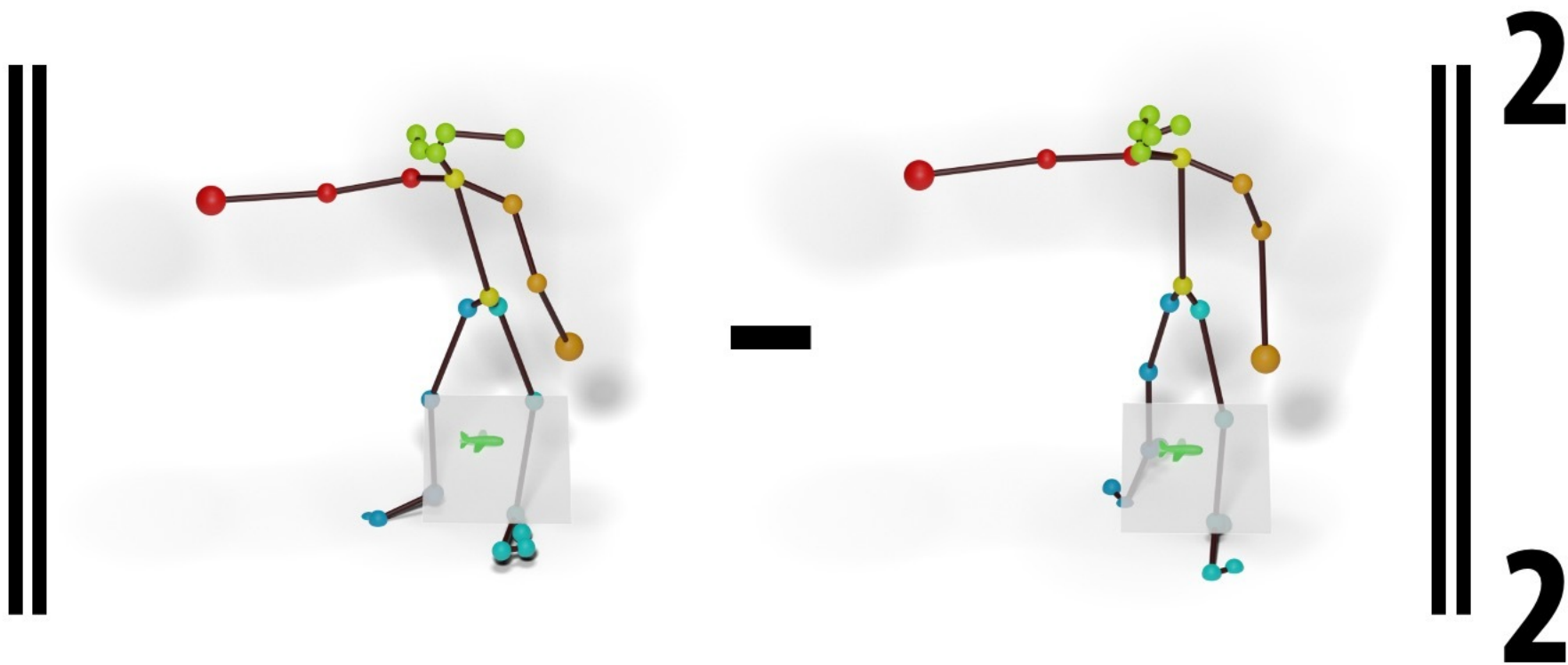| Input | Learning Trajectory Dependencies [1] | History Repeats Itself [2] | Ours | Target | |
|-------|--------------------------------------|----------------------------|------|--------|---|
| | | | | | "Pour" |
| | | | | | "Squeeze" |
| | | | | | "Pass" |

[1] Mao, Wei, et al. "Learning trajectory dependencies for human motion prediction." Proceedings of the IEEE/CVF International Conference on Computer Vision. 2019.

[2] Mao, Wei, Miaomiao Liu, and Mathieu Salzmann. "History repeats itself: Human motion prediction via motion attention." European Conference on Computer Vision. Springer, Cham, 2020.
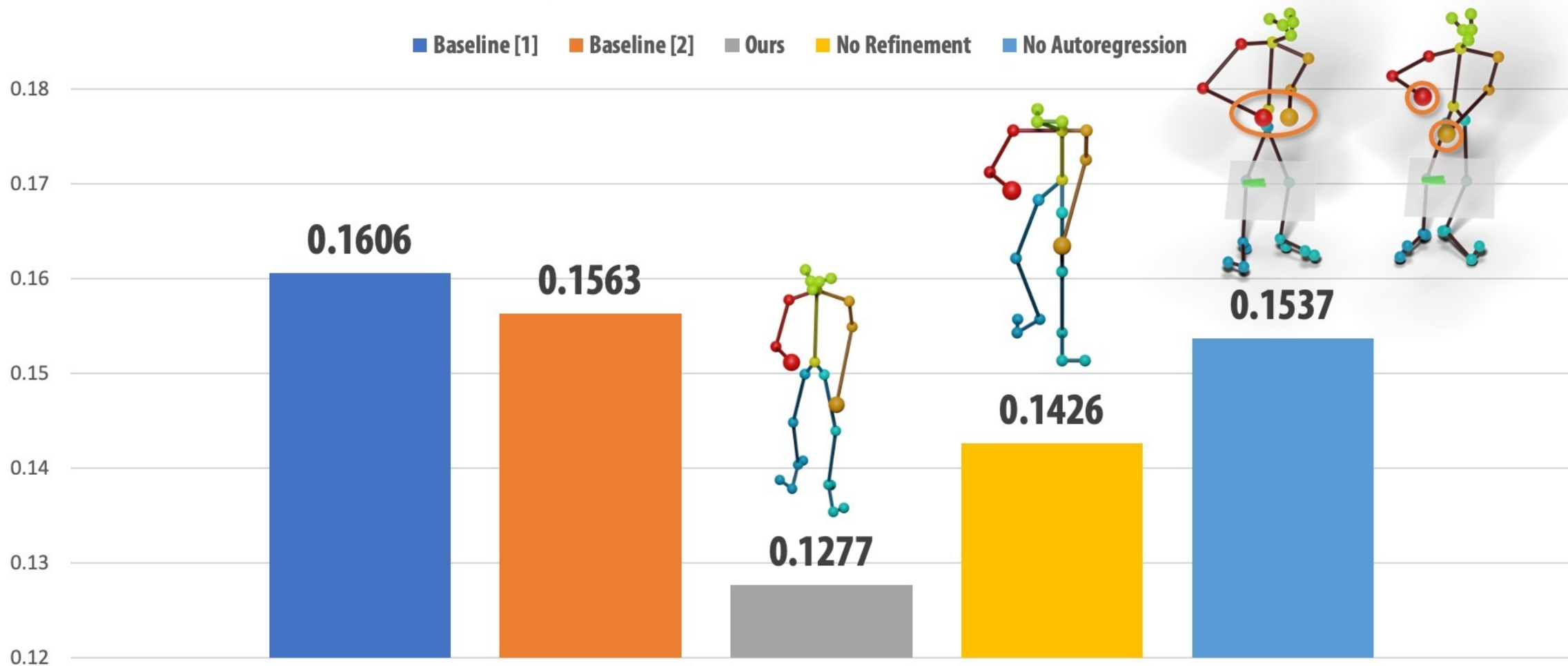
# Results: Mean Per-Joint Position Error

$$E_{\text{MPJPE}} = \frac{1}{25} \sum_{j=1}^{25} ||p_j' - p_j||_2^2$$
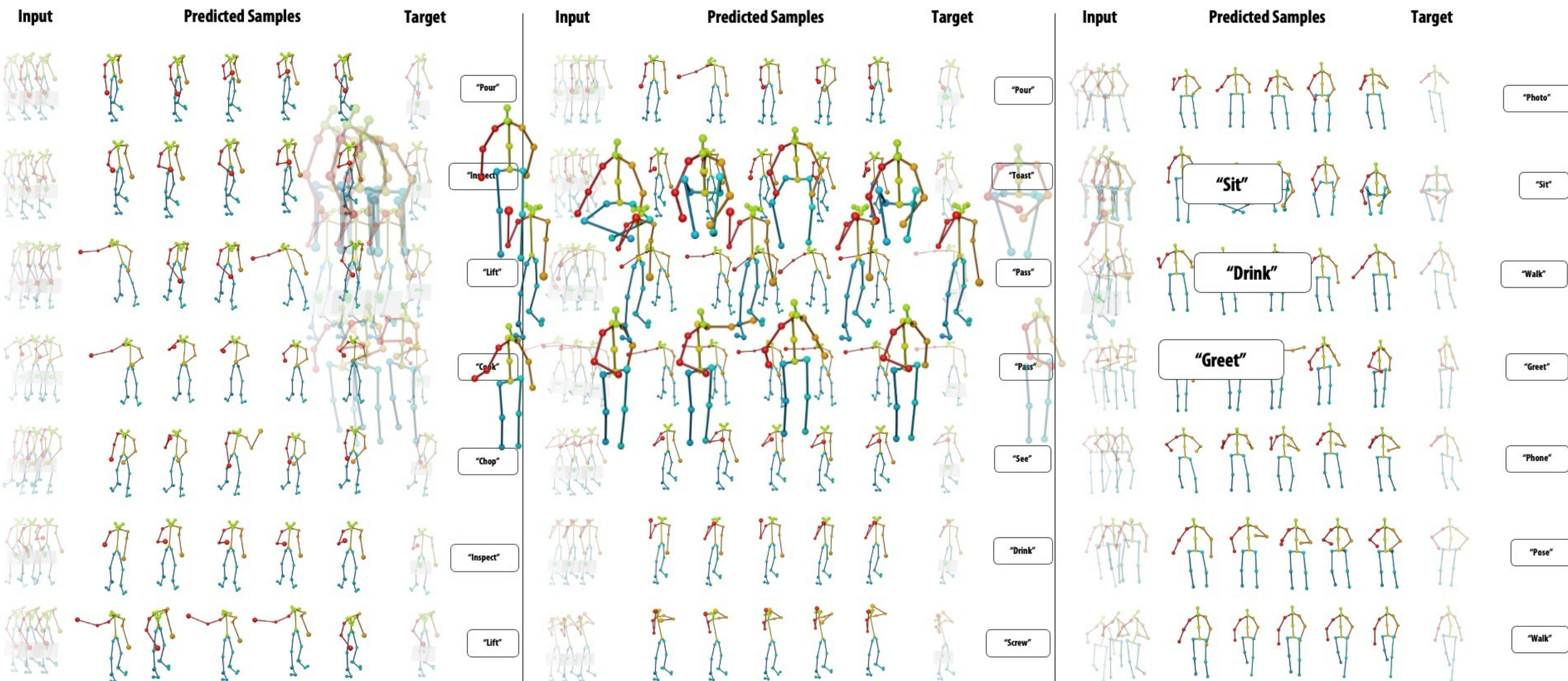
# Results: Quantitative



MPJPE = Mean Per-Joint Position Error

Baseline [1]  Baseline [2]  Ours  No Refinement  No Autoregression

| | | | | |
|---|---|---|---|---|
| 0.1606 | 0.1563 | 0.1277 | 0.1426 | 0.1537 |

[1] Mao, Wei, et al. "Learning trajectory dependencies for human motion prediction." Proceedings of the IEEE/CVF International Conference on Computer Vision. 2019.
[2] Mao, Wei, Miaomiao Liu, and Mathieu Salzmann. "History repeats itself: Human motion prediction via motion attention." European Conference on Computer Vision. Springer, Cham, 2020.

# Results: Qualitative – Multi-Modal Predictions

# 3D Human Behavior Generation: Action & Interaction

## Efficient Action Representation



**Forecasting Characteristic 3D Poses [1]**

## Complex Action Sequences



**FutureHuman3D [2]**

## Human-Object Interactions



**CG-HOI [3]**

[1] Diller, Christian, Thomas Funkhouser, and Angela Dai. "Forecasting characteristic 3d poses of human actions." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022.

[2] Diller, Christian, Thomas Funkhouser, and Angela Dai. "FutureHuman3D: Forecasting Complex Long-Term 3D Human Behavior from Video Observations." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024.

[3] Diller, Christian, and Angela Dai. "Cg-hoi: Contact-guided 3d human-object interaction generation." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024.

# 3D Human Behavior Generation: Action & Interaction
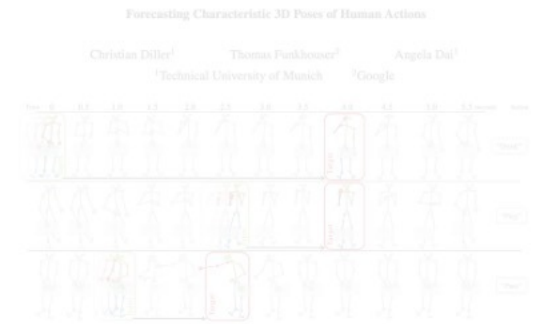
**Efficient Action Representation**     **Complex Action Sequences**     **Human-Object Interactions**



**Forecasting Characteristic 3D Poses [1]**     **FutureHuman3D [2]**     **CG-HOI [3]**

[1] Diller, Christian, Thomas Funkhouser, and Angela Dai. "Forecasting characteristic 3d poses of human actions." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022.

[2] Diller, Christian, Thomas Funkhouser, and Angela Dai. "FutureHuman3D: Forecasting Complex Long-Term 3D Human Behavior from Video Observations." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024.

[3] Diller, Christian, and Angela Dai. "Cg-hoi: Contact-guided 3d human-object interaction generation." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024.

# Related Work: Action Forecasting



**Past Observation Sequence**

**Forecasted Future**

# Related Work: 3D Pose Forecasting



$x_0$     $x_1$     $x_2$     $x_3$     $x_4$

3D Pose Forecasting

$y_0$     $y_1$     $y_2$     $y_3$

**Past Observation Sequence**

**Forecasted Future**

Time

# Task: Future Actions & 3D Poses from 2D

**2D RGB Images + Action Labels**

**3D Pose Sequence + Action Labels**



"take"  "cut"  "cut"

**Joint 3D Pose & Action Forecasting**

"screw"  "pour"  "screw"  "add"

**Past Observation Sequence**

**Forecasted Future**

25  Time

# Data: Uncorrelated 2D and 3D Human Poses

## 2D Action Sequences



4x

- Take
- Wash
- Take
- Take
- Take
- Close
- Take

- Take
- Peel
- Throw in Garbage
- Cut
- Add
- Throw in Garbage

## 3D Pose Data



**AMASS [1]**

**GRAB [2]**

**Human3.6m [3]**



**No correspondence**

[1] Mahmood, Naureen, et al. "AMASS: Archive of motion capture as surface shapes." Proceedings of the IEEE/CVF international conference on computer vision. 2019.
[3] Ionescu, Catalin, et al. "Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments." IEEE transactions on pattern analysis and machine intelligence 36.7 (2013): 1325-1339.
[2] Taheri, Omid, et al. "GRAB: A dataset of whole-body human grasping of objects." European conference on computer vision. Springer, Cham, 2020.

# Method: Architecture



Input Sequence

Forecasted 3D Pose + Action

2D Human Pose Observations

Pose History Encoder

"add"

Action CE Loss

Actions

"take"  "cut"  "change temperature"

Action Encoder

Characteristic 3D Pose

2D Proj.

Objects

drawer, frying pan, hand   knife, onion   hand, stove

Object Encoder

MLP Decoder

3D Adversarial Loss

2D Joint Loss

Time

# Results: Qualitative 3D Pose & Action – Cooking



**Input** / **Target**

"open"  "wash"  "dry"  "throw in garbage"  "take"    "wash"  "wash"  "shake"  "take"  "add"

**DLow** / **GSPS**

**STARS** / **Ours**

"wash"  "wash"  "move lid"  "shake"  "add"

# Results: Qualitative 3D Pose & Action – Furniture Assembly



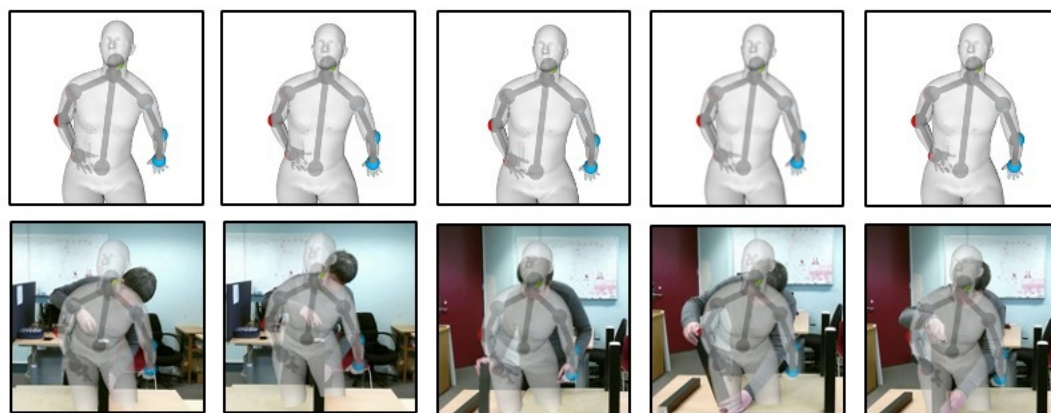Input: "rotate" "pick up" "align" "spin" "pick up"
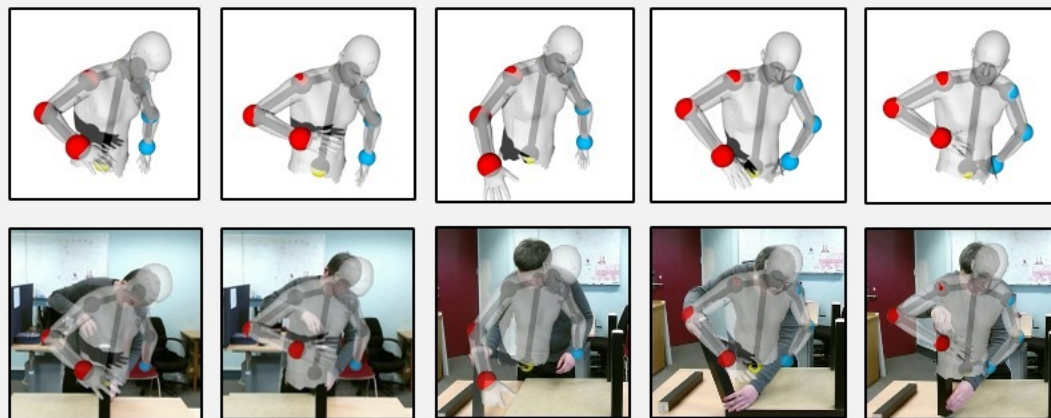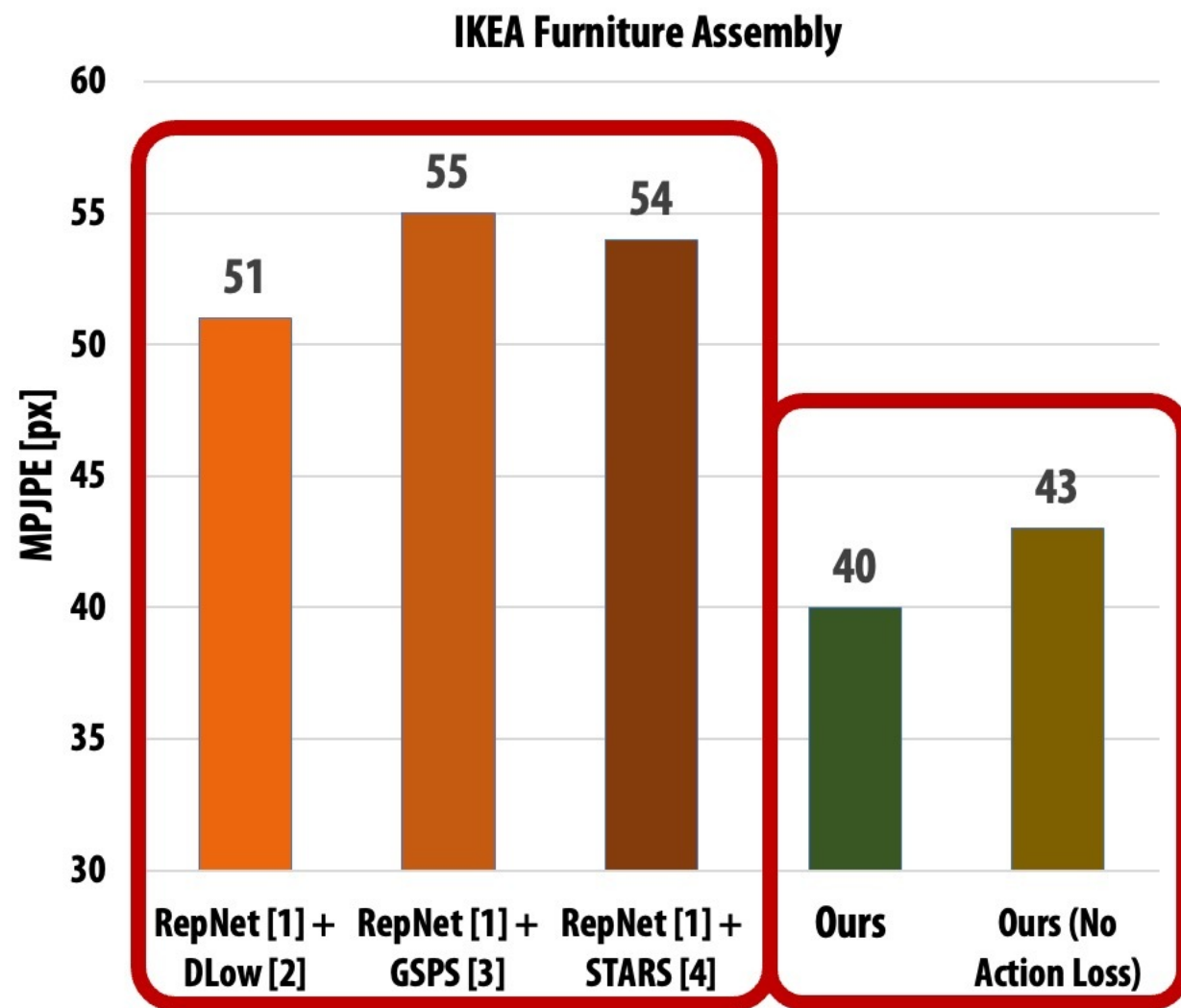Target: "align" "spin" "pick up" "align" "spin"

DLow
GSPS
STARS
Ours: "align" "spin" "align" "spin" "spin"

# Results: 3D Pose Forecasting – 2D Joint Error



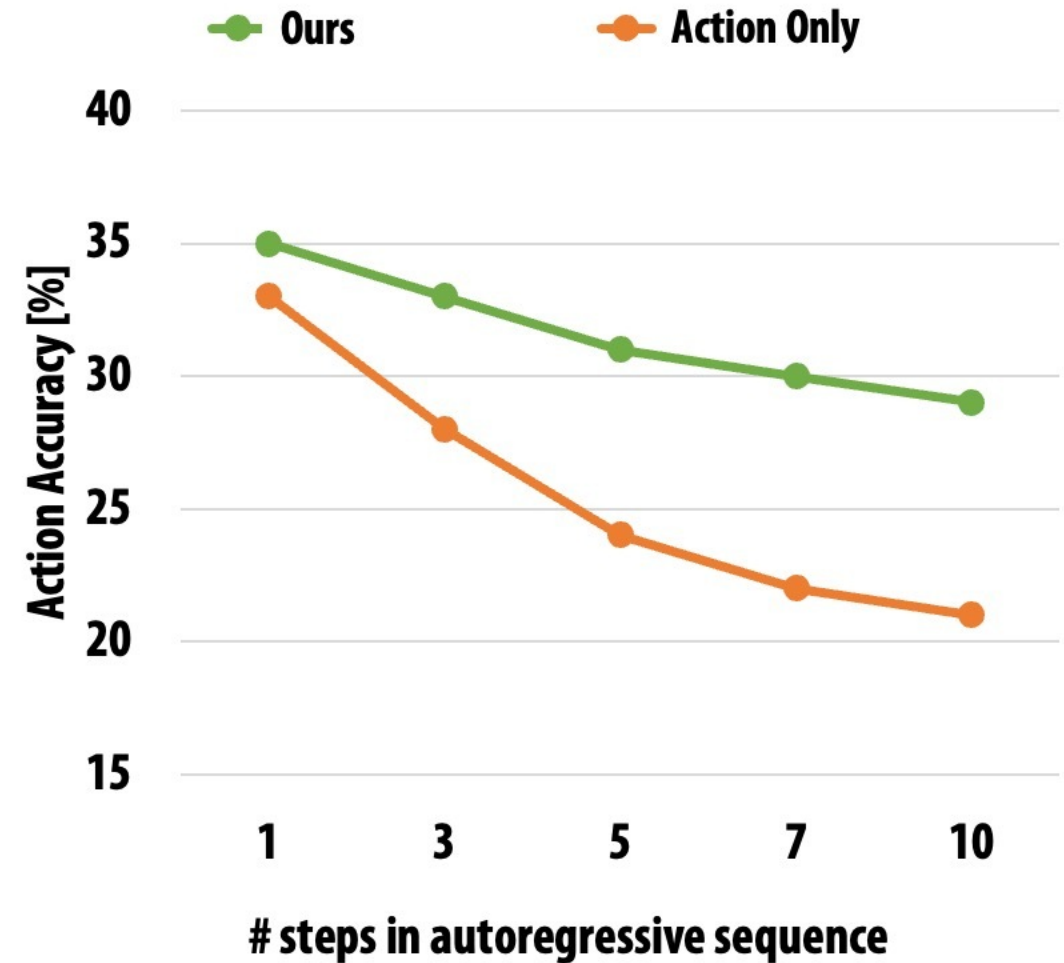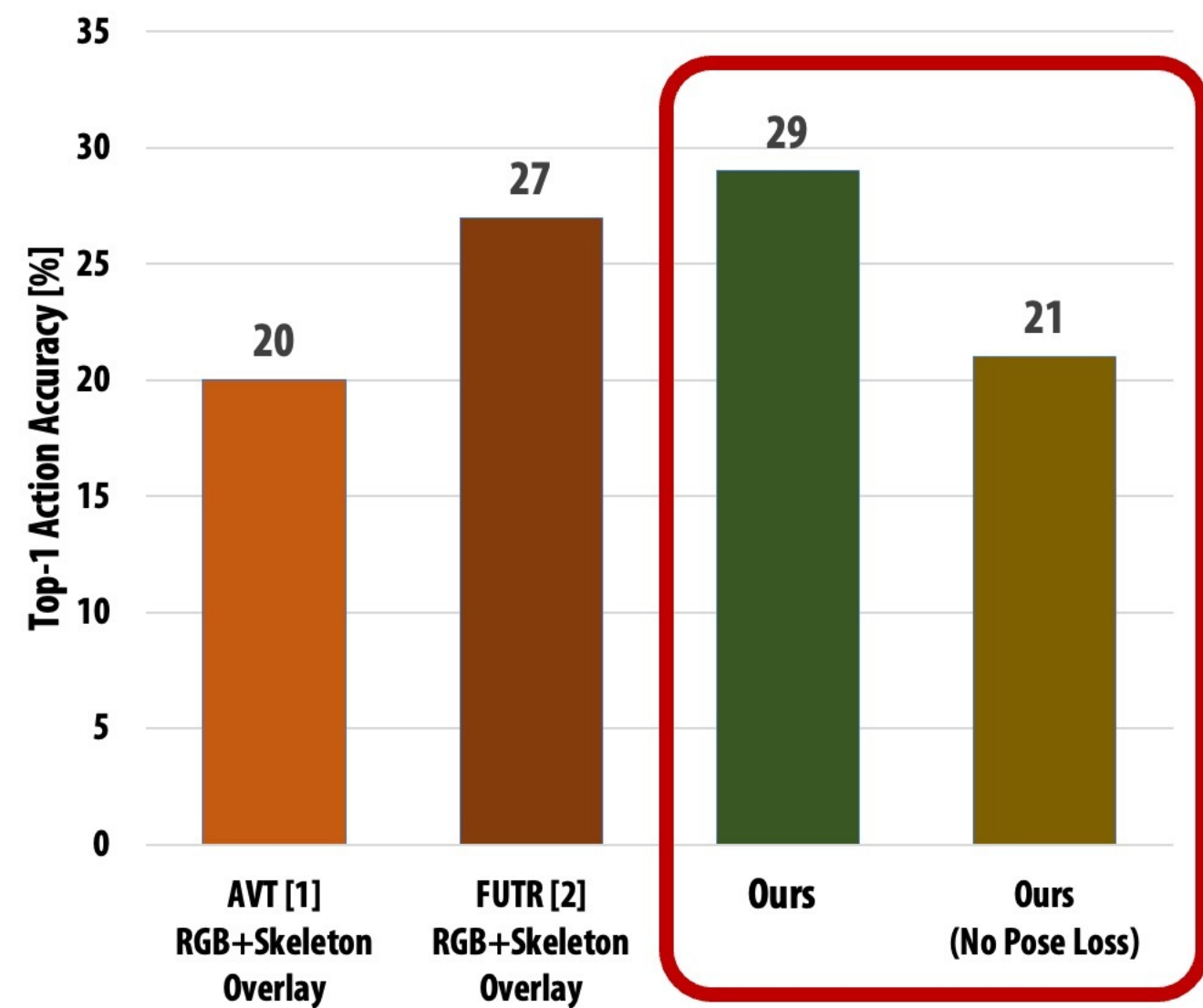**Cooking Sequences**

**IKEA Furniture Assembly**

[1] Wandt, Bastian, and Bodo Rosenhahn. "Repnet: Weakly supervised training of an adversarial reprojection network for 3d human pose estimation." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2019.
[2] Yuan, Ye, and Kris Kitani. "Dlow: Diversifying latent flows for diverse human motion prediction." ECCV 2020.
[3] Mao, Wei, Miaomiao Liu, and Mathieu Salzmann. "Generating Smooth Pose Sequences for Diverse Human Motion Prediction." Proceedings of the IEEE/CVF International Conference on Computer Vision. 2021.
[4] Xu, Sirui, Yu-Xiong Wang, and Liang-Yan Gui. "Diverse Human Motion Prediction Guided by Multi-level Spatial-Temporal Anchors." ECCV 2022.

# Results: Action Forecasting – Action Accuracy

[1] Girdhar, Rohit, and Kristen Grauman. "Anticipative video transformer." Proceedings of the IEEE/CVF International Conference on Computer Vision. 2021.
[2] Gong, Dayoung, et al. "Future transformer for long-term action anticipation." Proceedings of the IEEE/CVF International Conference on Computer Vision. 2022.

# 3D Human Behavior Generation: Action & Interaction
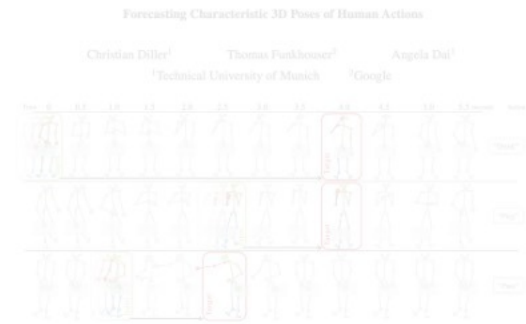
**Efficient Action Representation**

**Complex Action Sequences**

**Human-Object Interactions**



**Forecasting Characteristic 3D Poses [1]**

**FutureHuman3D [2]**

**CG-HOI [3]**

[1] Diller, Christian, Thomas Funkhouser, and Angela Dai. "Forecasting characteristic 3d poses of human actions." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022.

[2] Diller, Christian, Thomas Funkhouser, and Angela Dai. "FutureHuman3D: Forecasting Complex Long-Term 3D Human Behavior from Video Observations." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024.

[3] Diller, Christian, and Angela Dai. "Cg-hoi: Contact-guided 3d human-object interaction generation." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024.

# 3D Human Behavior Generation: Action & Interaction
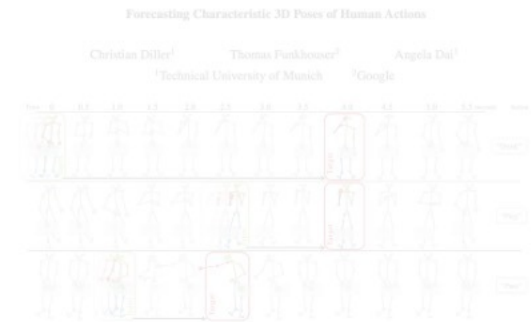
**Efficient Action Representation**

**Complex Action Sequences**

**Human-Object Interactions**



Forecasting Characteristic 3D Poses [1]

FutureHuman3D [2]

CG-HOI [3]

[1] Diller, Christian, Thomas Funkhouser, and Angela Dai. "Forecasting characteristic 3d poses of human actions." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022.

[2] Diller, Christian, Thomas Funkhouser, and Angela Dai. "FutureHuman3D: Forecasting Complex Long-Term 3D Human Behavior from Video Observations." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024.

[3] Diller, Christian, and Angela Dai. "Cg-hoi: Contact-guided 3d human-object interaction generation." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024.

# Task: Joint Human-Object Motion Generation

Move a chair with the hand

# Approach: Contact Modeling

# Method: Joint Training

# Method: Inference



Denoising U-Net

$D_H$

$D_O$

$D_C$

Recomputed Contact

Predicted Contact

38

# Results: Qualitative



Carry a suitcase

Move the chair backwards

Condition

Condition

# Results: Qualitative Comparison to Baseline MDM [1]



**Condition**

**MDM [1]**

**Ours**

[1] Guy Tevet, Sigal Raab, Brian Gordon, Yoni Shafir, Daniel Cohen-or, Amit Haim Bermano. "Human Motion Diffusion Model." The Eleventh International Conference on Learning Representations . 2023.

# Results: Qualitative Comparison to Baseline InterDiff [1]



Move a chair with the hand

Condition

InterDiff [1]

Ours

Carry a trash bin

[1] Xu, Sirui, et al. "Interdiff: Generating 3d human-object interactions with physics-informed diffusion." Proceedings of the IEEE/CVF International Conference on Computer Vision. 2023.

# Results: Ablation Study



No Contact Modeling      No Contact Guidance      Ours (Full)

# Results: Quantitative – User Study



**More Realistic**

MDM: Ours 81.8%, Baseline 18.2%
InterDiff: Ours 72.8%, Baseline 27.2%

**Follows Text Better**

MDM: Ours 79.5%, Baseline 20.5%
InterDiff: Ours 73.1%, Baseline 26.9%

Ours | Baseline

# Results: Quantitative

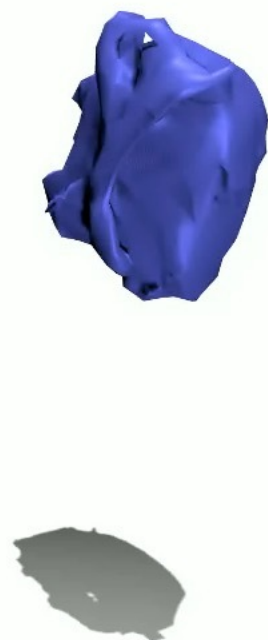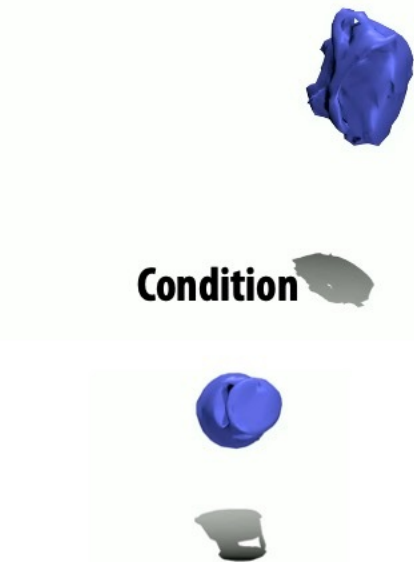| Task | Approach | BEHAVE | | | | CHAIRS | | | |
|------|----------|--------|--|--|--|--------|--|--|--|
| | | R-Prec. (top-3) ↑ | FID ↓ | Diversity → | MModality → | R-Prec. (top-3) ↑ | FID ↓ | Diversity → | MModality → |
| | Real (human) | 0.73 | 0.09 | 4.23 | 4.55 | 0.83 | 0.01 | 7.34 | 3.00 |
| Text-Cond. Human Only | MDM [71] | 0.52 | 4.54 | 5.44 | 5.12 | 0.72 | 5.99 | 6.83 | 3.45 |
| | InterDiff [84] | 0.49 | 5.36 | 3.98 | 3.98 | 0.63 | 6.76 | 5.24 | 2.44 |
| | **Ours** | **0.60** | **4.26** | **4.92** | **4.10** | **0.78** | **5.24** | **7.90** | **3.22** |
| | Real | 0.81 | 0.17 | 6.80 | 6.24 | 0.87 | 0.02 | 9.91 | 6.12 |
| Motion-Cond. HOI | InterDiff [84] | 0.68 | 3.86 | 5.62 | 5.90 | 0.67 | 4.83 | 7.49 | 4.87 |
| | **Ours** | **0.71** | **3.52** | **6.89** | **6.43** | **0.79** | **4.01** | **8.42** | **6.29** |
| Text-Cond. HOI | MDM [71] | 0.49 | 9.21 | 6.51 | 8.19 | 0.53 | 9.23 | 6.23 | 7.44 |
| | InterDiff [84] | 0.53 | 8.70 | 3.85 | 4.23 | 0.69 | 7.53 | 5.23 | 4.63 |
| | **Ours** | **0.62** | **6.31** | **6.63** | **5.47** | **0.74** | **6.45** | **8.91** | **5.94** |

# Application: Object Trajectory Guidance



Carry a backpack on the back

Condition

Move a stool

Generation

Condition

Generation

# Application: 3D Static Scene Population

# 3D Human Behavior Generation: Action & Interaction



Efficient Action Representation

Forecasting Characteristic 3D Poses [1]

Complex Action Sequences

FutureHuman3D [2]

**Human-Object Interactions**

**CG-HOI [3]**

[1] Diller, Christian, Thomas Funkhouser, and Angela Dai. "Forecasting characteristic 3d poses of human actions." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022.

[2] Diller, Christian, Thomas Funkhouser, and Angela Dai. "FutureHuman3D: Forecasting Complex Long-Term 3D Human Behavior from Video Observations." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024.

[3] Diller, Christian, and Angela Dai. "Cg-hoi: Contact-guided 3d human-object interaction generation." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024.

# 3D Human Behavior Generation: Action & Interaction
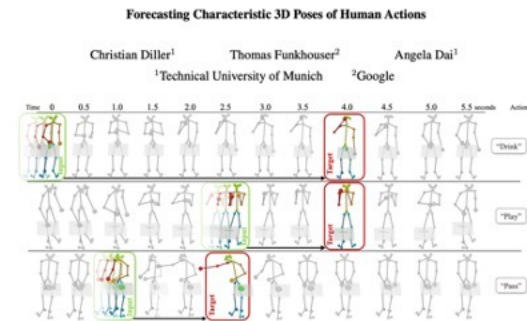
**Efficient Action Representation**



**Forecasting Characteristic 3D Poses [1]**

**Complex Action Sequences**



**FutureHuman3D [2]**

**Human-Object Interactions**



**CG-HOI [3]**

[1] Diller, Christian, Thomas Funkhouser, and Angela Dai. "Forecasting characteristic 3d poses of human actions." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022.

[2] Diller, Christian, Thomas Funkhouser, and Angela Dai. "FutureHuman3D: Forecasting Complex Long-Term 3D Human Behavior from Video Observations." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024.

[3] Diller, Christian, and Angela Dai. "Cg-hoi: Contact-guided 3d human-object interaction generation." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024.

# Summary and Conclusion



"Drink"   "Pour"

- **Predicting future characteristic 3D poses of human actions**
  - **Probabilistic approach for capturing the most likely future 3D action poses**

"screw"   "pour"

- **Forecasting complex long-term 3D human behavior from 2D**
  - **Joint action and 3D pose forecasting of composite long-term behavior**

Play with a yoga ball   Lift a small table

- **Contact-Guided 3D Human-Object Interactions**
  - **Realistic human-object interaction generation from text and geometry**

# Outlook: 3D Scene Understanding

**Reconstruction [1]**

**Semantic Instance Segmentation [2]**

**Affordance Prediction [3]**

[1] Dai, Angela, Christian Diller and Matthias Nießner. "SG-NN: Sparse Generative Neural Networks for Self-Supervised Scene Completion of RGB-D Scans." 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2020): 846-855.

[2] Hou, Ji, Angela Dai and Matthias Nießner. "3D-SIS: 3D Semantic Instance Segmentation of RGB-D Scans." 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2019): 4416-4425.

[3] Savva, Manolis, Angel X. Chang, Pat Hanrahan, Matthew Fisher, and Matthias Nießner. "SceneGrok: Inferring action maps in 3D environments." ACM transactions on graphics (TOG) 33, no. 6 (2014): 1-10.

51

# Outlook: Dynamic Human Interactions in 3D Scenes



**Text-based motion and interaction [1]**

The person walks forward from the curtain to pick up his guitar.

The person cartwheels towards the campfire from the table.

**Zero-shot path-finding with large language models [2]**

[1] Yi, Hongwei, et al. "Generating human interaction motions in scenes with text control." European Conference on Computer Vision. Springer, Cham, 2024.

[2] Qu, Haoxuan, Ziyan Guo, and Jun Liu. "GPT-Connect: Interaction between Text-Driven Human Motion Generator and 3D Scenes in a Training-free Manner." arXiv preprint arXiv:2403.14947 (2024).

# Thank You!

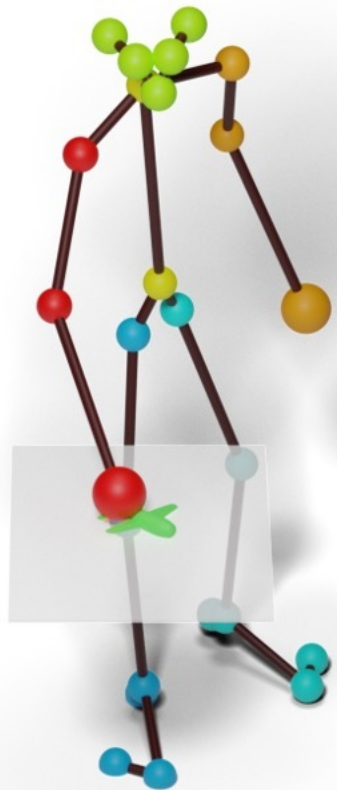Prof. Angela Dai          Prof. Michael Black          Prof. Stefan Leutenegger          Prof. Thomas Funkhouser

# 3D Human Behavior Generation through Action and Interaction Synthesis

**Christian Diller**

**Supervisor: Prof. Angela Dai**

**Tuesday, 10th December 2024**