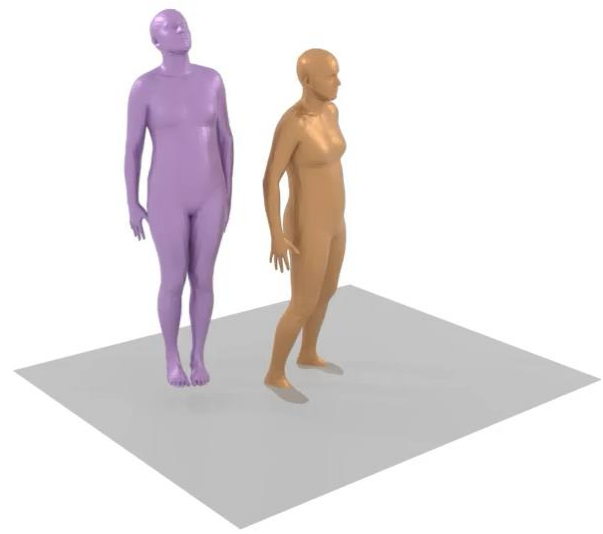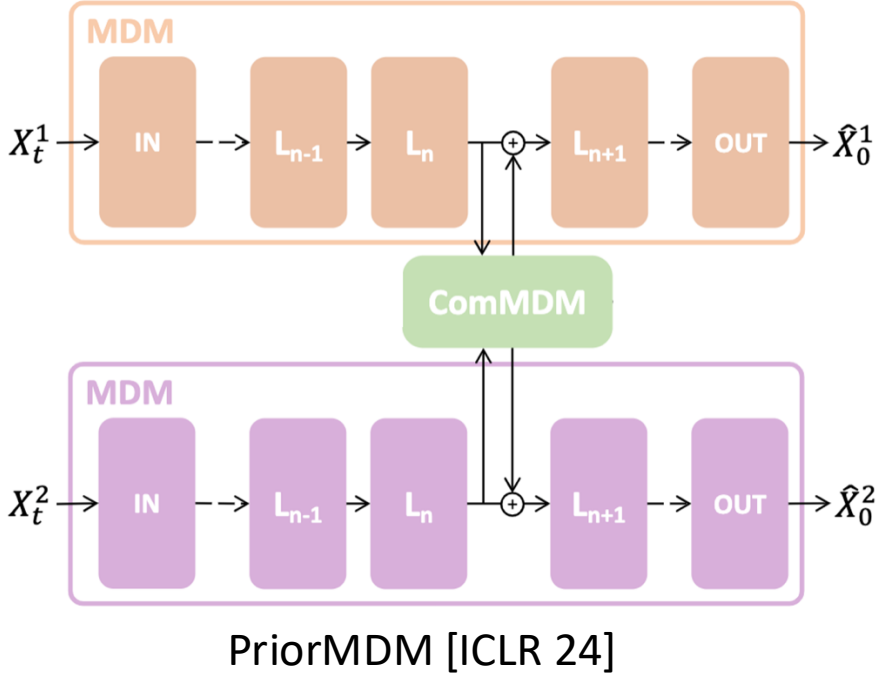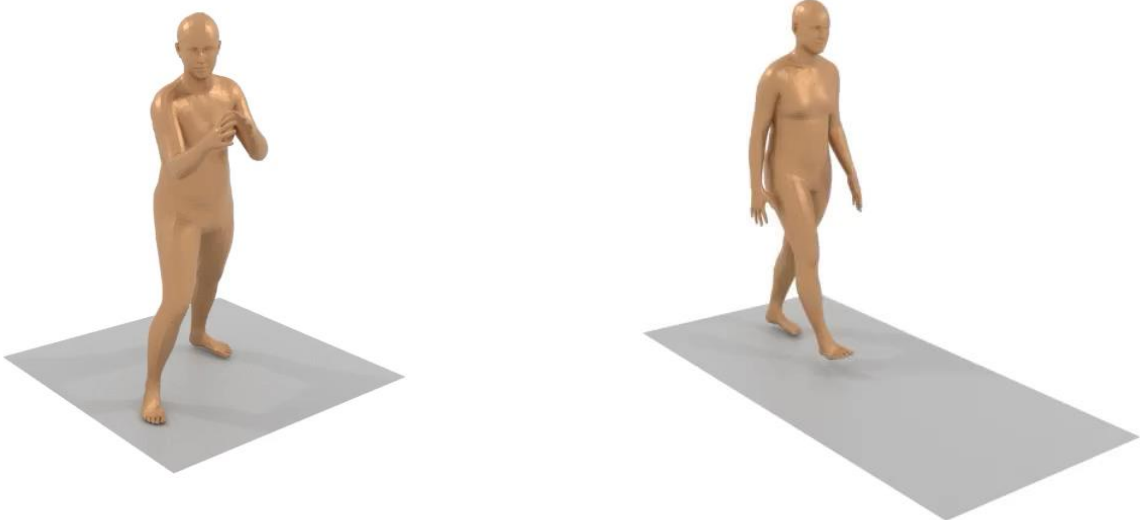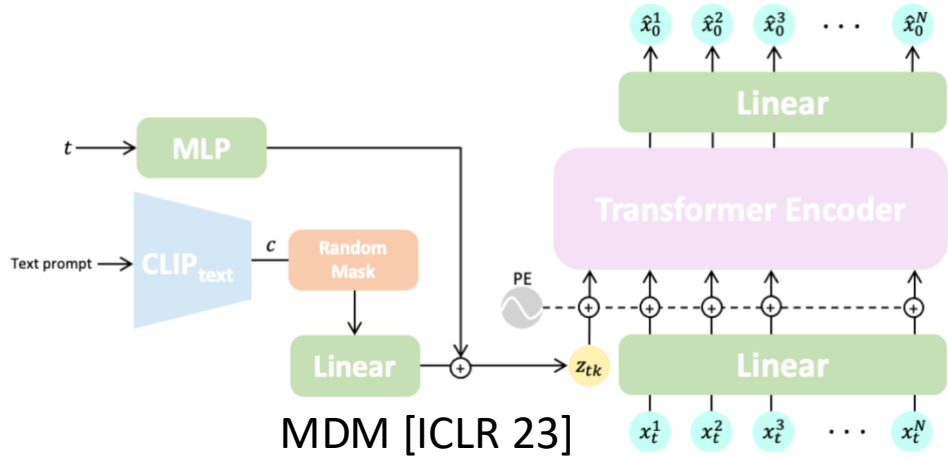# InterControl: Zero-shot Human Interaction Generation by Controlling Every Joint

## NeurIPS 2024

Zhenzhi Wang, Jingbo Wang, Yixuan Li, Dahua Lin, Bo Dai
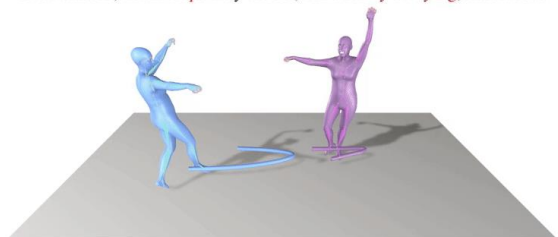
CUHK, Shanghai AI Lab, HKU

Code: https://github.com/zhenzhiwang/intercontrol

# Human Motion Generation



MDM [ICLR 23]

PriorMDM [ICLR 24]

# Multi-Person Interaction



InterHuman [IJCV24]

Inter-X [CVPR 24]

Our work was done in Nov, 2023. Thus, back to that time, could we generate multi-person interactions only by single-person motion data?

Advantages: the interaction will not be restricted by only two people defined by the dataset.

# Definition of Interaction



Instruction: I would like to play video games for a while. After that, I will go to sleep.

UniHSI [ICLR24]

Yes! We find that interaction could be represented by joint distances. And the semantics of interactions could be understood by Large Language Models.

# Spatial Control

# Zero-shot Interaction Generation

# Our interaction are realistic



Three-legged race   Hold both hands   Hold hands   Hold hands in dancing   Hold hands in groups

**(c) Physics animation with human-wise interactions**

**(a) Daily life**

Reference Interactions

Generated by ours

Kick foot with foot   Kick head with foot   Hold hands   Hit head in 2v1 fighting

**(b) Fighting**

Reference Interactions

Generated by ours

# Quantitative comparisons

Table 1: **Spatial control** results on HumanML3D [14]. → means closer to real data is better. *Random One/Two/Three* reports the average performance over 1/2/3 randomly selected joints in evaluation. [†] means our evaluation on their model.

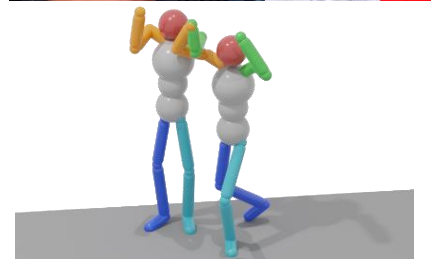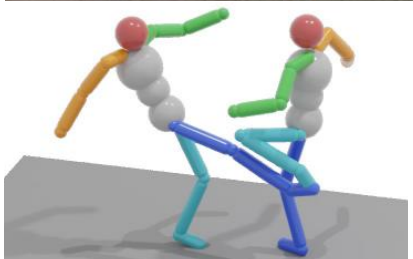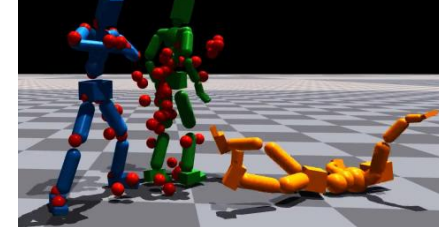| Method | Joint | FID ↓ | R-precision ↑ (Top-3) | Diversity → | Foot skating ratio ↓ | Traj. err. ↓ (50 cm) | Loc. err. ↓ (50 cm) | Avg. err. ↓ (m) |
|---|---|---|---|---|---|---|---|---|
| Real data | - | 0.002 | 0.797 | 9.503 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| MDM [55] | No Control | 0.544 | 0.611 | 9.446 | 0.0943 | 0.8909 | 0.6015 | 1.1843 |
| PriorMDM [51][†] | | 0.498 | 0.586 | 9.167 | 0.0924 | 0.3726 | 0.2210 | 0.4552 |
| GMD [27][†] | Root | 0.276 | 0.655 | 9.245 | 0.1108 | 0.0987 | 0.0356 | 0.1457 |
| OmniControl [65] | | 0.218 | 0.687 | 9.422 | **0.0547** | 0.0387 | 0.0096 | **0.0338** |
| Ours | | **0.159** | 0.671 | 9.482 | 0.0729 | **0.0132** | **0.0004** | 0.0496 |
| OmniControl [65] | Random one | 0.310 | **0.693** | **9.502** | 0.0608 | 0.0617 | 0.0107 | 0.0404 |
| Ours | | 0.178 | 0.669 | 9.498 | 0.0968 | 0.0403 | 0.0031 | 0.0741 |
| Ours | Random two | 0.184 | 0.670 | 9.410 | 0.0948 | 0.0475 | 0.0030 | 0.0911 |
| Ours | Random three | 0.199 | 0.673 | 9.352 | 0.0930 | 0.0487 | 0.0026 | 0.0969 |

Table 2: Evaluation on (left) spatial errors and (right) user preference in interactions.

| Spatial Errors | Traj. err. (20 cm) ↓ | Loc. err. (20 cm) ↓ | Avg. err. (m) ↓ |
|---|---|---|---|
| PriorMDM [51] | 0.6931 | 0.3487 | 0.6723 |
| Ours | **0.0082** | **0.0005** | **0.0084** |

| User-study | Preference |
|---|---|
| PriorMDM [51] | 18.8% |
| Ours | **81.2%** |

# Quantitative comparisons

Table 3: **Ablation studies** on the HumanML3D [14] dataset.

| Item | Method | FID ↓ | R-precision ↑ (Top-3) | Diversity → | Foot skating ratio ↓ | Traj. err. ↓ (50 cm) | Loc. err. ↓ (50 cm) | Avg. err.↓ (m) |
|------|--------|-------|-----------------------|-------------|----------------------|----------------------|---------------------|----------------|
| (1) | Ours (random joint) | **0.178** | 0.669 | **9.498** | 0.0968 | 0.0403 | 0.0031 | 0.0741 |
| (2) | w/o ControlNet | 0.965 | 0.621 | 9.216 | 0.1624 | 0.0879 | 0.0059 | 0.1013 |
| (3) | w/ original $c$ | 0.227 | 0.656 | 9.544 | 0.1004 | 0.0697 | 0.0042 | 0.0785 |
| (4) | w/o IK guidance | 0.187 | 0.664 | 9.598 | **0.0704** | 0.8569 | 0.4553 | 0.6557 |
| (5) | IK guidance on $x_0$ | 0.211 | 0.668 | 9.394 | 0.1164 | 0.0907 | 0.0088 | 0.0981 |
| (6) | w/ 1-st order grad | 0.198 | 0.668 | 9.472 | 0.0987 | 0.0879 | 0.0096 | 0.0877 |
| (7) | sparsity = 0.25 | 0.248 | **0.671** | 9.442 | 0.0801 | 0.0106 | 0.0007 | 0.0546 |
| (8) | sparsity = 0.025 | 0.255 | 0.663 | 9.520 | 0.0705 | **0.0015** | **0.0001** | **0.0067** |

Table 4: **Inference time analysis** on a NVIDIA A100 GPU.

| Sub-Modules | MDM | + Control Module | + Guidance $t \in [10, 999]$ | + Guidance $t \in [0, 9]$ |
|-------------|-----|------------------|------------------------------|---------------------------|
| Time (s) | 39.1 | 57.3 | 76.5 | 80.1 |

# Qualitative comparisons



Description: One person leaps to execute a kick towards the other, and the other jumps to respond with a kick.

PriorMDM

Ours

Figure 3: Comparison with PriorMDM [51] in **user-study** of zero-shot human interaction generation.



Interaction Description

Two people are dancing with each other, and they hold hands from time to time.

Two people are fighting with the third one by using their hands and foots.
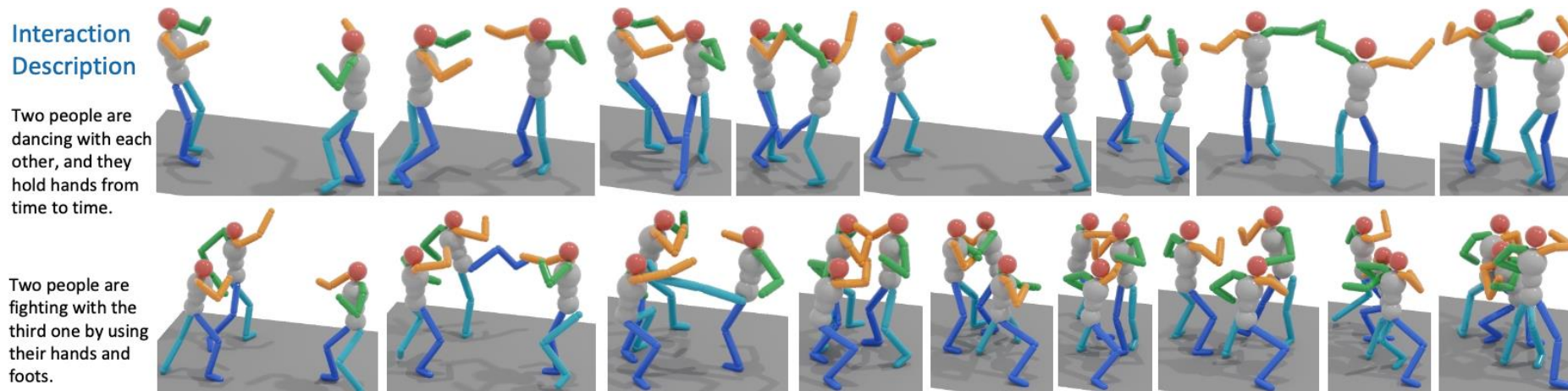
Figure 4: **Qualitative results** of zero-shot human interaction generation.

# Thank you!