

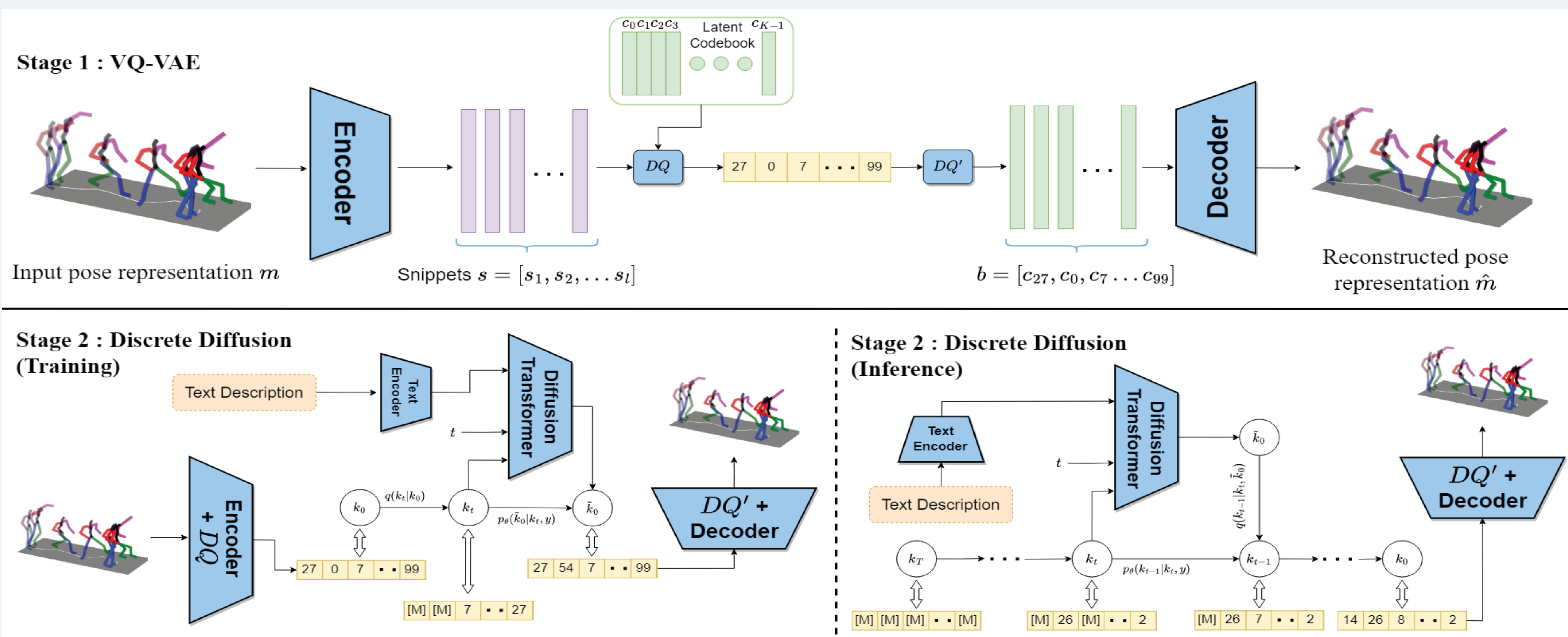
## Motivation

- Synthesizing well aligned human motions based on input conditions is an essential task with many applications in gaming, simulation, and virtual reality.
- For a many-to-many task like text-motion synthesis, probabilistic methods such as Diffusion are necessary. Despite, excellent generative abilities, diffusion is expensive on raw motions and so faster alternatives are needed.
- Hence, we use latent space Discrete Diffusion with the assumption that human motion can be efficiently represented by discrete sequence of small motion snippets.

## Contributions

- We model the text-to-motion generation task as a discrete denoising diffusion probabilistic model, which allows reduced diffusion steps for faster inferences while producing high quality results.
- Evaluated our method (MoDDM) in a comparison with state-of-the-art methods using both objective metrics and subjective user study. The results demonstrated that our method outperforms the previous methods in both motion quality and text-to-motion matching accuracy.

## Two-stage Architecture consisting of VQ-VAE and Discrete Diffusion Model

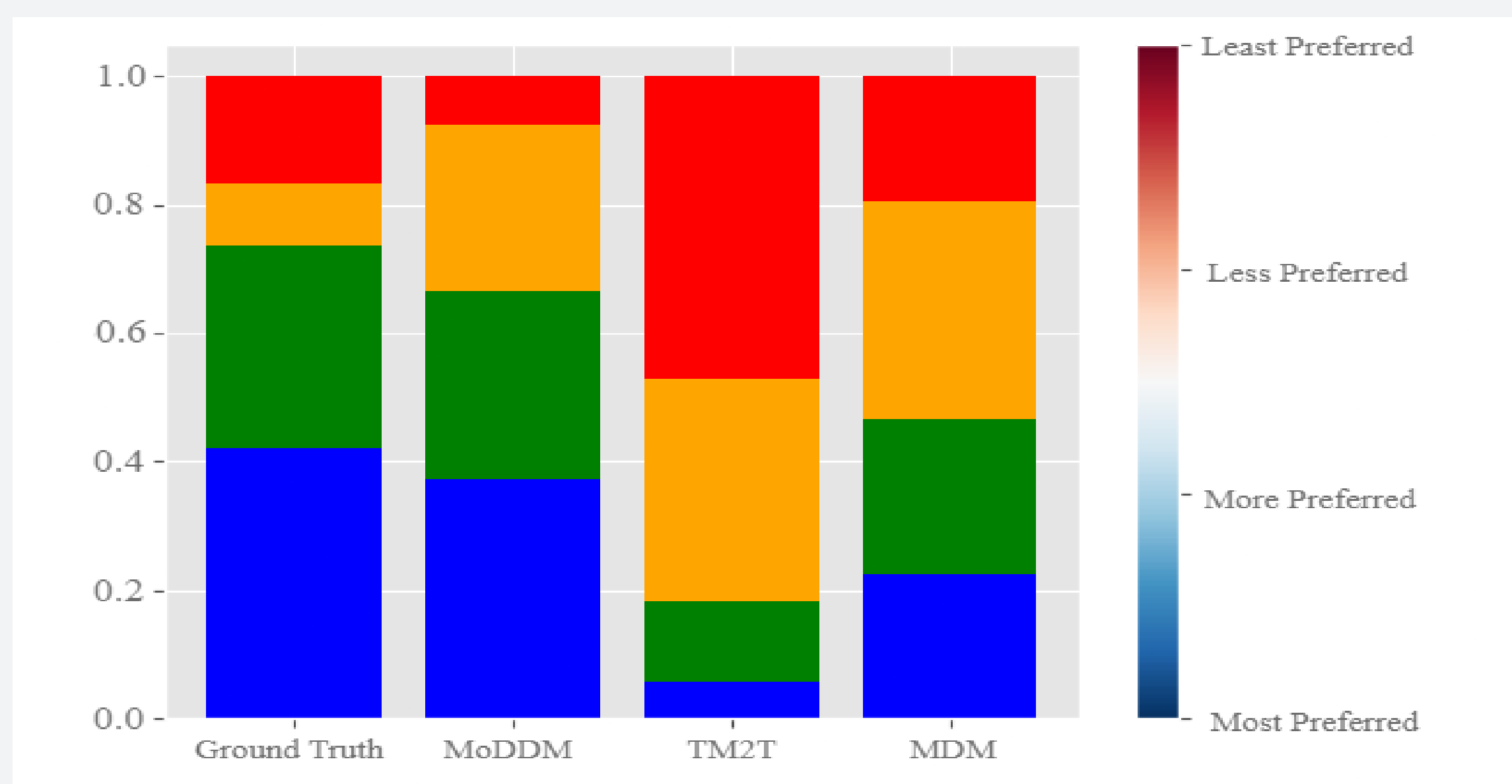


## Quantitative Evaluation on HumanML3D Test Set

Methods	R Precision $\uparrow$			FID $\downarrow$	MultiModal Dist $\downarrow$	Diversity $\rightarrow$
	Top 1	Top 2	Top 3			
<b>Real Motions</b>	0.511 $\pm$ .003	0.703 $\pm$ .003	0.797 $\pm$ .002	0.002 $\pm$ .000	2.974 $\pm$ .008	9.503 $\pm$ .065
Seq2Seq [1]	0.180 $\pm$ .002	0.300 $\pm$ .002	0.396 $\pm$ .002	11.75 $\pm$ .035	5.529 $\pm$ .007	6.223 $\pm$ .061
Language2Pose [2]	0.246 $\pm$ .002	0.387 $\pm$ .002	0.486 $\pm$ .002	11.02 $\pm$ .046	5.296 $\pm$ .008	7.676 $\pm$ .058
MDM [3]	-	-	0.611 $\pm$ .007	0.544 $\pm$ .044	5.566 $\pm$ .027	<b>9.559</b> $\pm$ .086
TM2T [4]	0.424 $\pm$ .003	<b>0.618</b> $\pm$ .003	<b>0.729</b> $\pm$ .002	1.501 $\pm$ .017	<b>3.467</b> $\pm$ .011	8.589 $\pm$ .076
MoDDM (Ours)	<b>0.425</b> $\pm$ .004	<u>0.615</u> $\pm$ .004	<u>0.713</u> $\pm$ .003	<b>0.294</b> $\pm$ .006	<u>3.553</u> $\pm$ .009	<u>9.178</u> $\pm$ .093

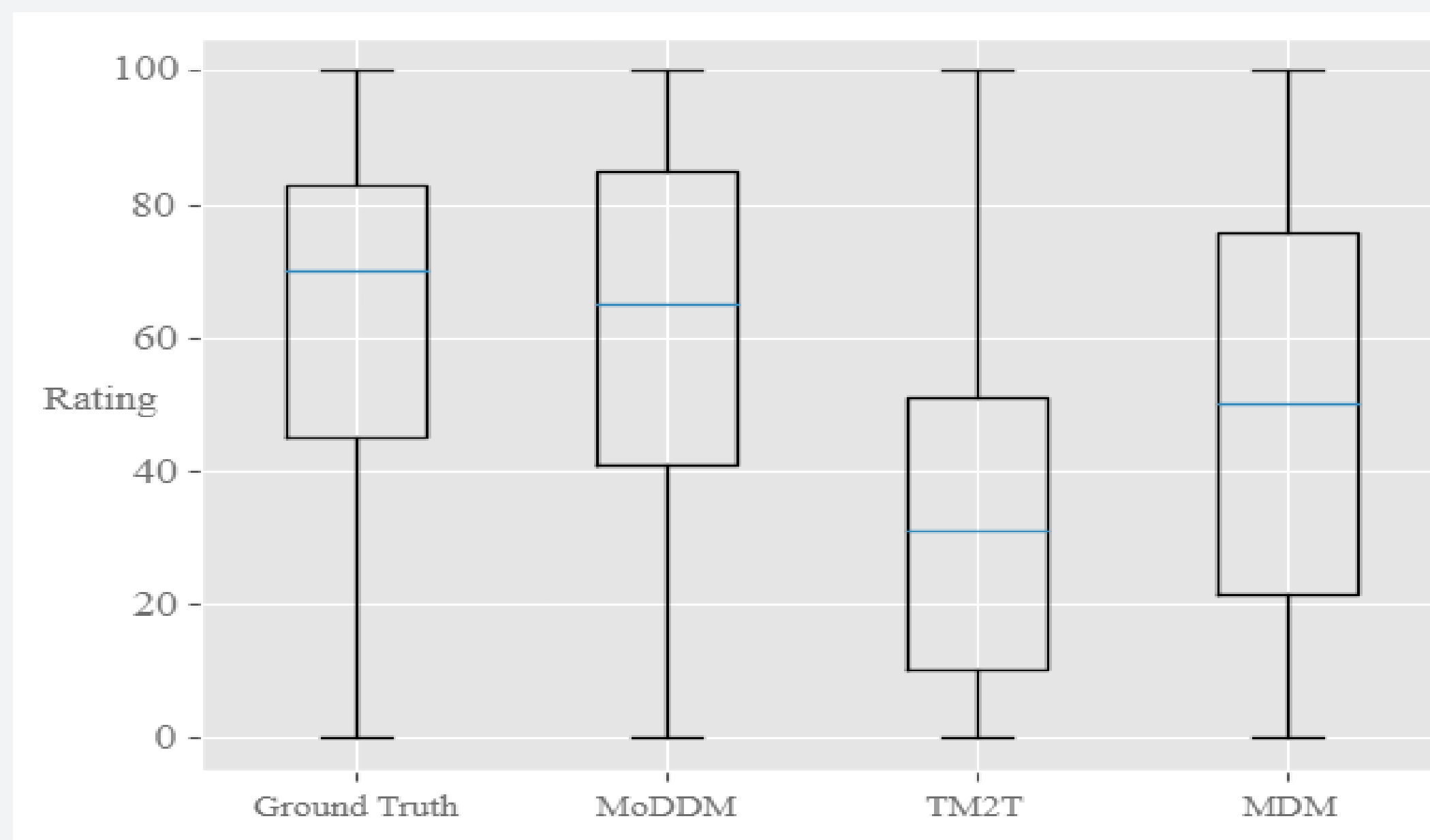
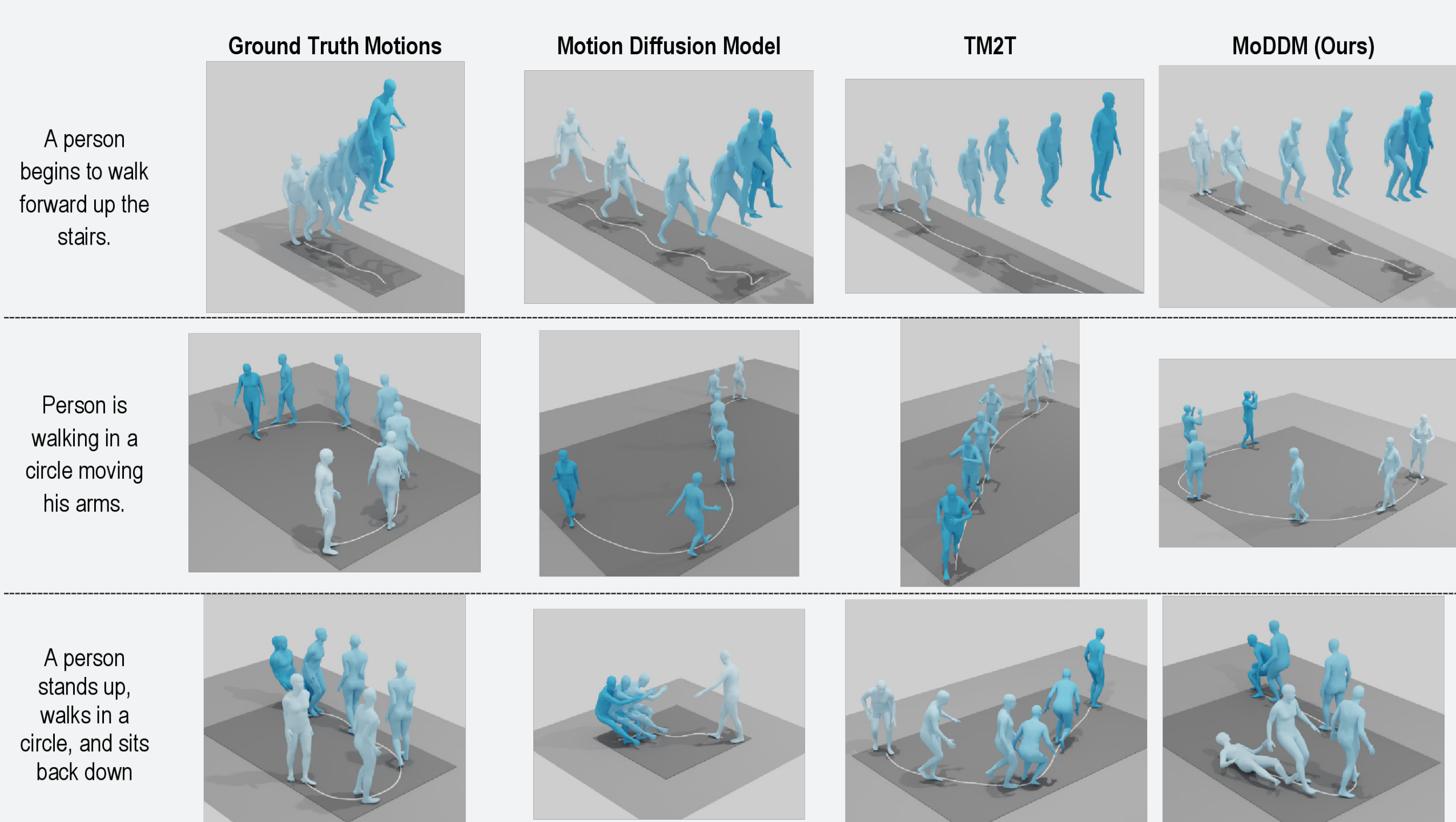
$\pm$  indicates 95% confidence interval, and  $\rightarrow$  means the closer to Real motions the better. Bold face indicates the best result, while underscore refers to the second best.

## Evaluation of Motion Alignment to Text



For ground truth and each comparison method, a color bar indicates the percentage of its preference levels.

## Qualitative Comparisons on HumanML3D Test Set



Boxes cover 25th and 75th percentiles, and whiskers represent the 5th and 95th percentiles. Box notches represent median values.