

Motion Mamba: Efficient and Long Sequence Motion Generation

ECCV 2024



Zeyu Zhang^{12*}†



Akide Liu^{1*}



Ian Reid³



Richard Hartley²



Bohan Zhuang¹



Hao Tang⁴✉

*Equal Contribution. ✉Corresponding author.

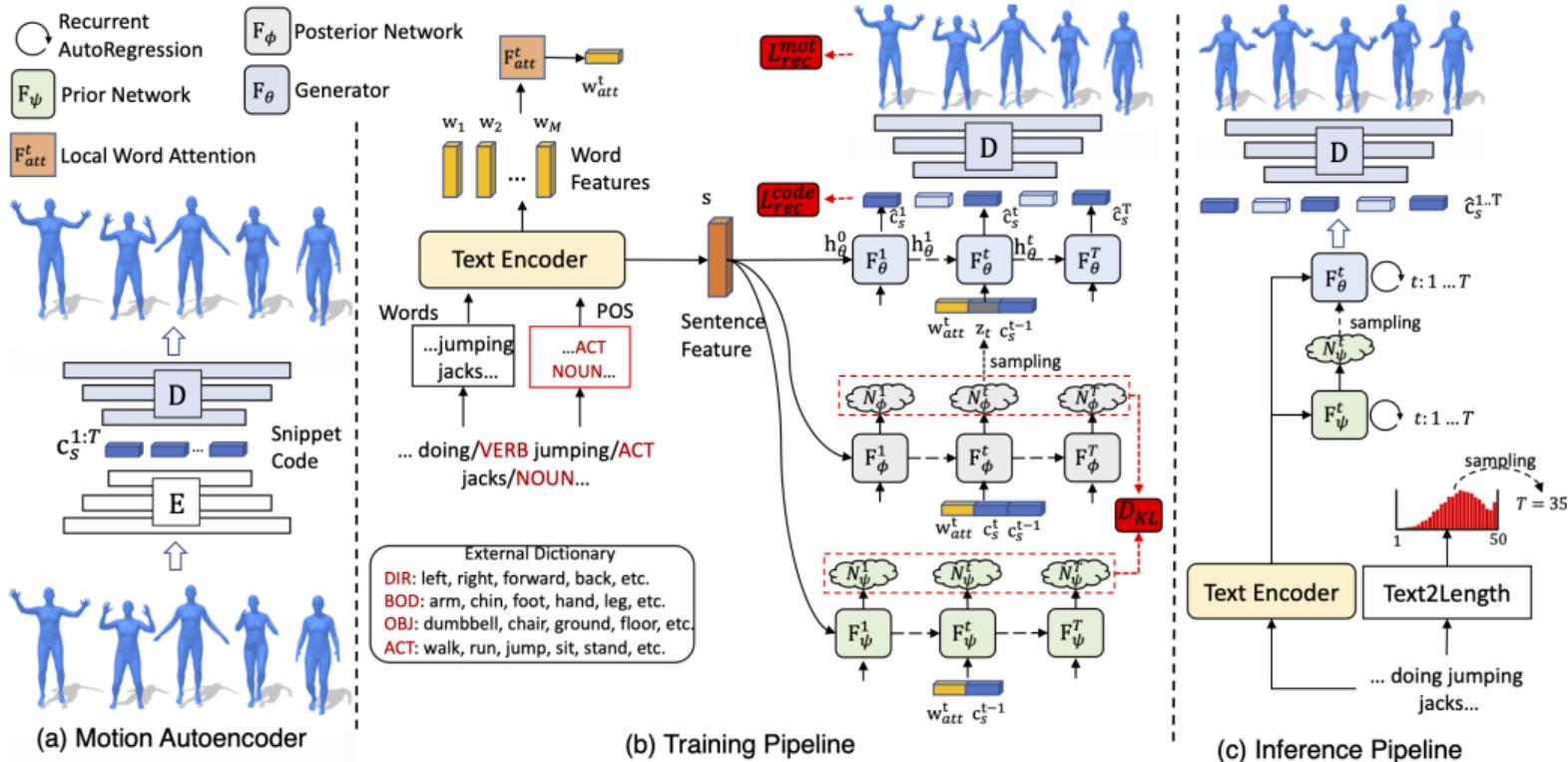
†Work done while being a research assistant at Monash University.

¹Monash University ²The Australian National University

³Mohamed bin Zayed University of Artificial Intelligence ⁴Peking University

Talk @ miHoYo Shanghai, July 22, 2024

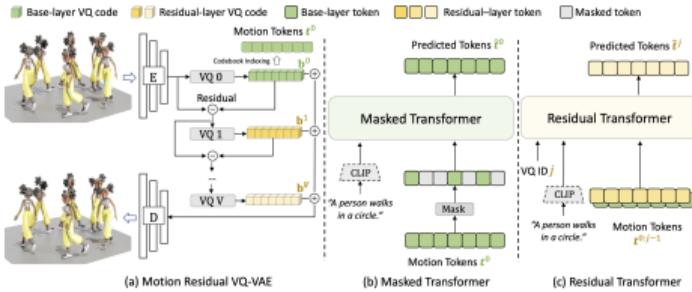
Text-to-Motion



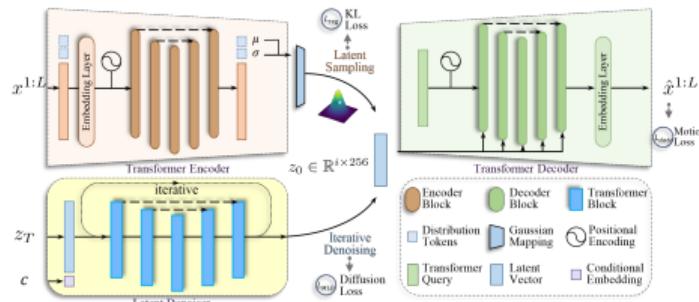
T2M for Now

- Autoregressive Method
 - Jiang et al. *MotionGPT: Human Motion as a Foreign Language* (NIPS 2023)
 - Pinyoanuntapong et al. *MMM: Generative Masked Motion Model* (CVPR 2024)
 - Guo et al. *MoMask: Generative Masked Modeling of 3D Human Motions* (CVPR 2024)
- Diffusion-based Method
 - Zhang et al. *MotionDiffuse: Text-Driven Human Motion Generation with Diffusion Model* (TPAMI 2024)
 - Tevet et al. *MDM: Human Motion Diffusion Model* (ICLR 2023)
 - Chen et al. *MLD: Executing your Commands via Motion Diffusion in Latent Space* (CVPR 2023)
 - Zhang et al. **Motion Mamba: Efficient and Long Sequence Motion Generation** (ECCV 2024)

MoMask vs MLD

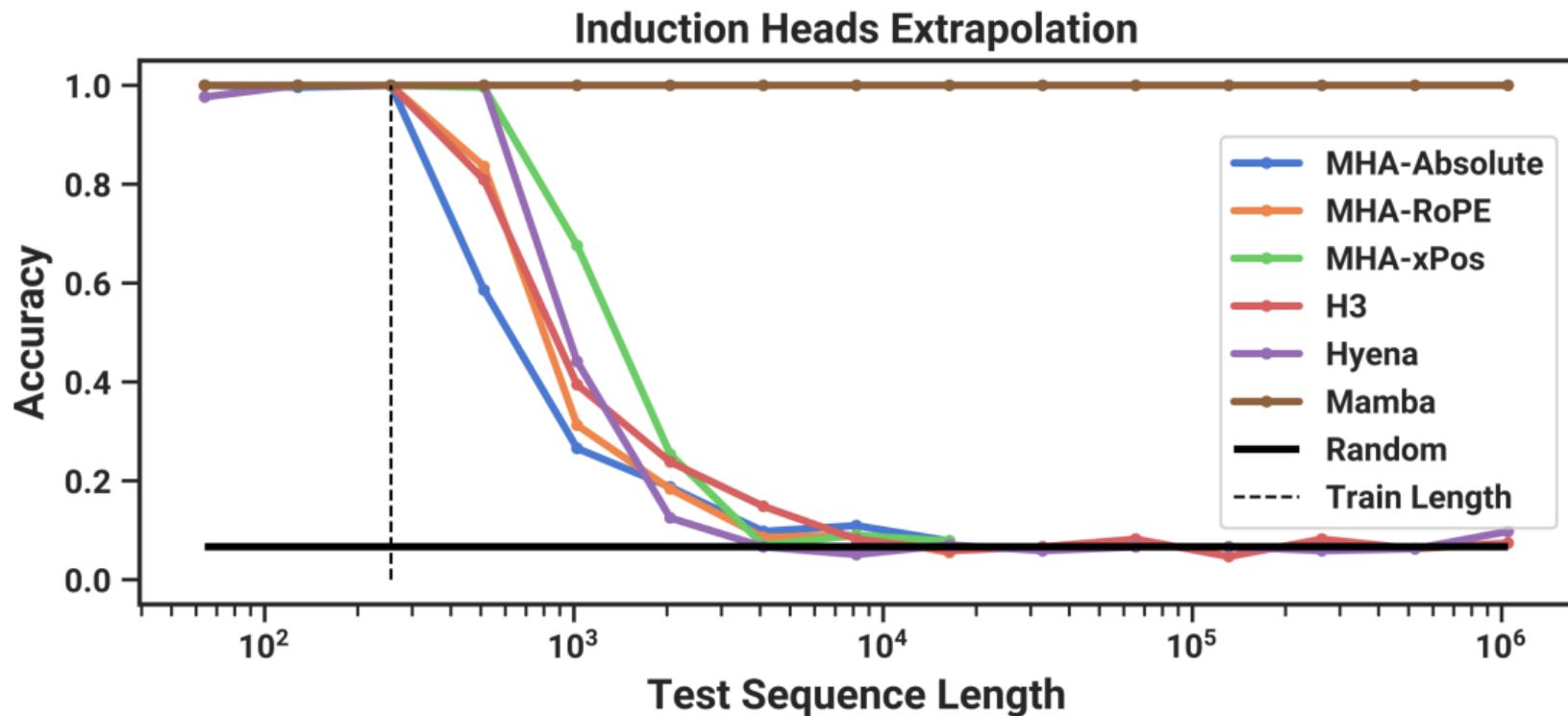


Guo et al. *MoMask: Generative Masked Modeling of 3D Human Motions* (CVPR 2024)

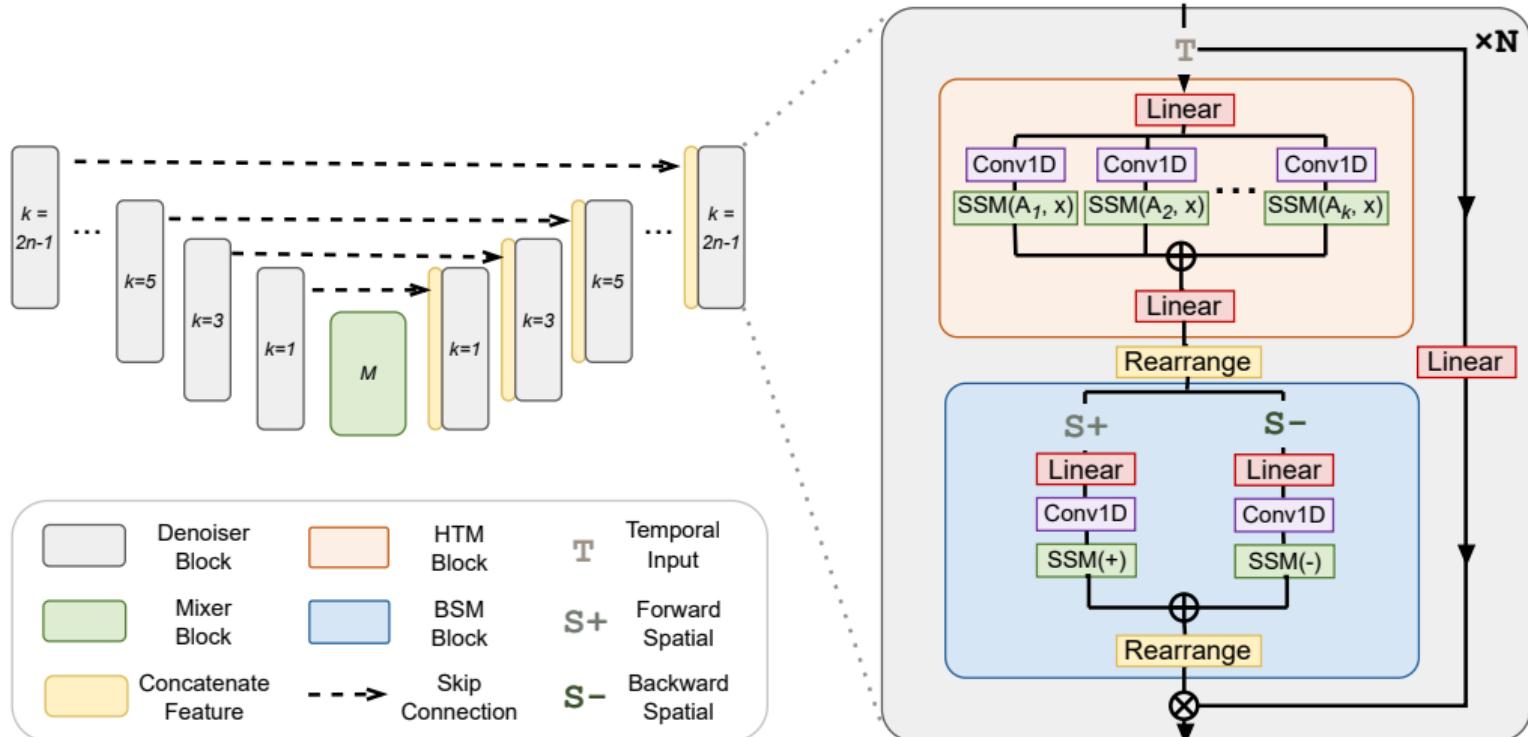


Chen et al. *MLD: Executing your Commands via Motion Diffusion in Latent Space* (CVPR 2023)

Motivation: Efficient and Long Sequence Modeling



Motion Mamba Architecture



Comparative Studies on HumanML3D

Comparison of text-conditional motion synthesis on **HumanML3D**. The motion encoder from MLD evaluates these metrics. The evaluation results are sorted by descending FIDs. The right arrow → refers to that closer to real motion. **Bold** and underline indicate the best and second best result.

Method	R Precision ↑			FID↓	MM Dist↓	Diversity→	MModality↑
	Top 1	Top 2	Top 3				
Real	0.511 ^{±.003}	0.703 ^{±.003}	0.797 ^{±.002}	0.002 ^{±.000}	2.974 ^{±.008}	9.503 ^{±.065}	-
Seq2Seq	0.180 ^{±.002}	0.300 ^{±.002}	0.396 ^{±.002}	11.75 ^{±.035}	5.529 ^{±.007}	6.223 ^{±.061}	-
LJ2P	0.246 ^{±.001}	0.387 ^{±.002}	0.486 ^{±.002}	11.02 ^{±.046}	5.296 ^{±.008}	7.676 ^{±.058}	-
T2G	0.165 ^{±.001}	0.267 ^{±.002}	0.345 ^{±.002}	7.664 ^{±.030}	6.030 ^{±.008}	6.409 ^{±.071}	-
Hier	0.301 ^{±.002}	0.425 ^{±.002}	0.552 ^{±.004}	6.532 ^{±.024}	5.012 ^{±.018}	8.332 ^{±.042}	-
TEMOS	0.424 ^{±.002}	0.612 ^{±.002}	0.722 ^{±.002}	3.734 ^{±.028}	3.703 ^{±.008}	8.973 ^{±.071}	0.368 ^{±.018}
T2M	0.457 ^{±.002}	0.639 ^{±.003}	0.740 ^{±.003}	1.067 ^{±.002}	3.340 ^{±.008}	9.188 ^{±.002}	2.090 ^{±.083}
MDM	0.320 ^{±.005}	0.498 ^{±.004}	0.611 ^{±.007}	0.544 ^{±.044}	5.566 ^{±.027}	9.559 ^{±.086}	2.799 ^{±.072}
MotionDiffuse	0.491 ^{±.001}	0.681 ^{±.001}	0.782 ^{±.001}	0.630 ^{±.001}	3.113 ^{±.001}	9.410 ^{±.049}	1.553 ^{±.042}
MLD	0.481 ^{±.003}	0.673 ^{±.003}	0.772 ^{±.002}	0.473 ^{±.013}	3.196 ^{±.010}	9.724 ^{±.082}	2.413 ^{±.079}
Motion Mamba (Ours)	0.502^{±.003}	0.693^{±.002}	0.792^{±.002}	0.281^{±.009}	3.060^{±.058}	9.871 ^{±.084}	2.294 ^{±.058}

Comparative Studies on KIT-ML

We involve the **KIT-ML** dataset and evaluate the SOTA methods on the text-to-motion task. The evaluation results are sorted by descending FIDs. The right arrow → refers to that closer to real motion. **Bold** and underline indicate the best and second best result.

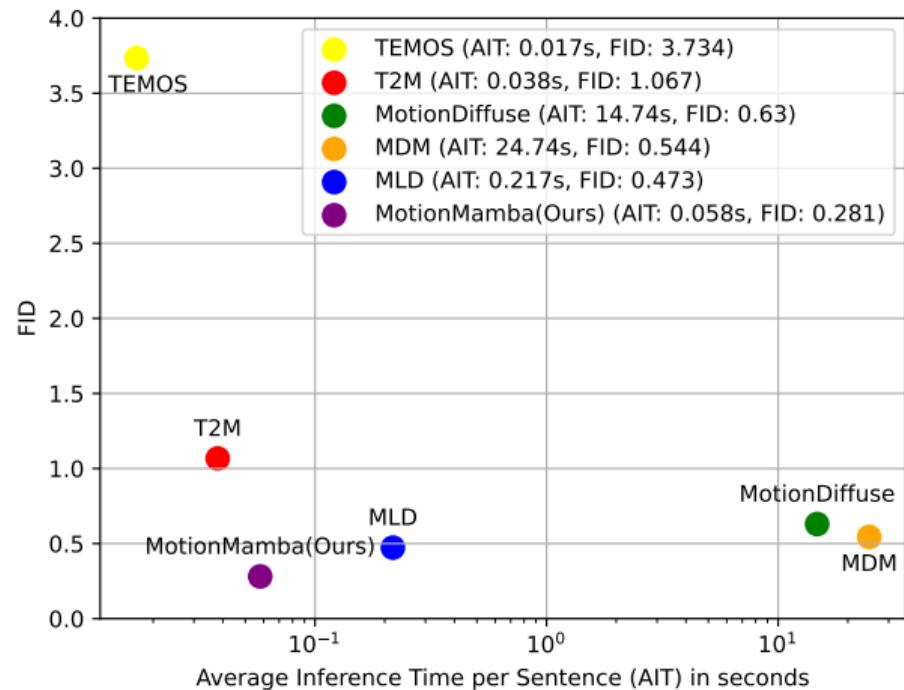
Method	R Precision ↑			FID↓	MM Dist↓	Diversity→	MModality↑
	Top 1	Top 2	Top 3				
Real	0.424 ^{±.005}	0.649 ^{±.006}	0.779 ^{±.006}	0.031 ^{±.004}	2.788 ^{±.012}	11.08 ^{±.097}	-
Seq2Seq	0.103 ^{±.003}	0.178 ^{±.005}	0.241 ^{±.006}	24.86 ^{±.348}	7.960 ^{±.031}	6.744 ^{±.106}	-
T2G	0.156 ^{±.004}	0.255 ^{±.004}	0.338 ^{±.005}	12.12 ^{±.183}	6.964 ^{±.029}	9.334 ^{±.079}	-
LJ2P	0.221 ^{±.005}	0.373 ^{±.004}	0.483 ^{±.005}	6.545 ^{±.072}	5.147 ^{±.030}	9.073 ^{±.100}	-
Hier	0.255 ^{±.006}	0.432 ^{±.007}	0.531 ^{±.007}	5.203 ^{±.107}	4.986 ^{±.027}	9.563 ^{±.072}	<u>2.090</u> ^{±.083}
TEMOS	0.353 ^{±.006}	0.561 ^{±.007}	0.687 ^{±.005}	3.717 ^{±.051}	3.417 ^{±.019}	10.84 ^{±.100}	0.532 ^{±.034}
T2M	0.370 ^{±.005}	0.569 ^{±.007}	0.693 ^{±.007}	2.770 ^{±.109}	3.401 ^{±.008}	10.91 ^{±.119}	1.482 ^{±.065}
MDM	0.164 ^{±.004}	0.291 ^{±.004}	0.396 ^{±.004}	0.497 ^{±.021}	9.191 ^{±.022}	10.85 ^{±.109}	1.907 ^{±.214}
MotionDiffuse	<u>0.417</u> ^{±.004}	<u>0.621</u> ^{±.004}	<u>0.739</u> ^{±.004}	1.954 ^{±.062}	2.958 ^{±.005}	11.10 ^{±.143}	0.730 ^{±.013}
MLD	0.390 ^{±.008}	0.609 ^{±.008}	0.734 ^{±.007}	<u>0.404</u> ^{±.027}	3.204 ^{±.027}	10.80 ^{±.117}	2.192 ^{±.071}
Motion Mamba (Ours)	0.419 ^{±.006}	0.645 ^{±.005}	0.765 ^{±.006}	0.307 ^{±.041}	<u>3.021</u> ^{±.025}	<u>11.02</u> ^{±.098}	1.678 ^{±.064}

Long Sequence Motion Generation

In order to evaluate the models' capability in long sequence motion generation, we compared our method with an existing approach on the recently introduced HumanML3D-LS dataset. This dataset comprises motion sequences longer than 190 frames from the original evaluation set. Our model demonstrates superior performance compared to other methods.

Method	R Precision ↑			FID↓	MM Dist↓	Diversity→	MModality↑
	Top 1	Top 2	Top 3				
Real	$0.437 \pm .003$	$0.622 \pm .004$	$0.721 \pm .004$	$0.004 \pm .000$	$3.343 \pm .015$	$8.423 \pm .090$	-
MDM	$0.368 \pm .005$	$0.553 \pm .006$	$0.672 \pm .005$	$0.802 \pm .044$	$3.860 \pm .025$	$8.817 \pm .068$	-
MotionDiffuse	$0.367 \pm .004$	$0.521 \pm .004$	$0.623 \pm .004$	$2.460 \pm .062$	$3.789 \pm .005$	$8.707 \pm .143$	$1.602 \pm .013$
MLD	$0.403 \pm .005$	$0.584 \pm .005$	$0.690 \pm .005$	$0.952 \pm .020$	$3.580 \pm .016$	$9.050 \pm .085$	$2.711 \pm .104$
Motion Mamba (Ours)	$0.417 \pm .003$	$0.606 \pm .003$	$0.713 \pm .004$	$0.668 \pm .019$	$3.435 \pm .015$	$9.021 \pm .070$	$2.373 \pm .084$

Efficiency



Visualization

MotionDiffuse



MDM



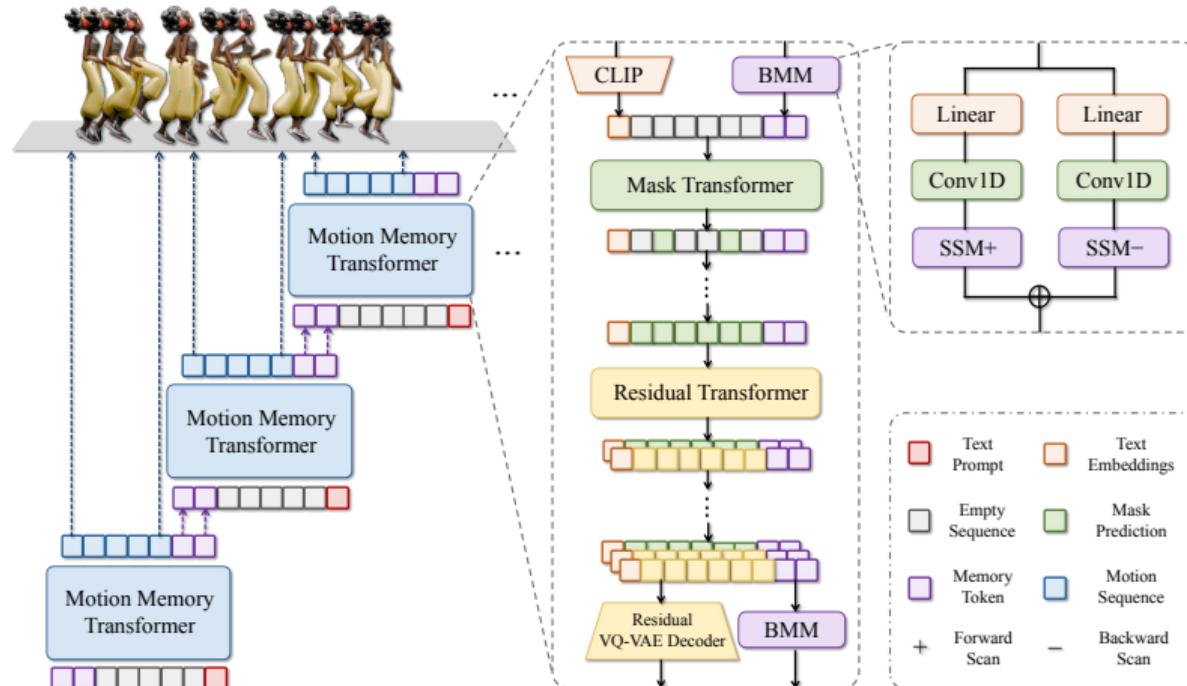
MLD



Motion Mamba (Ours)

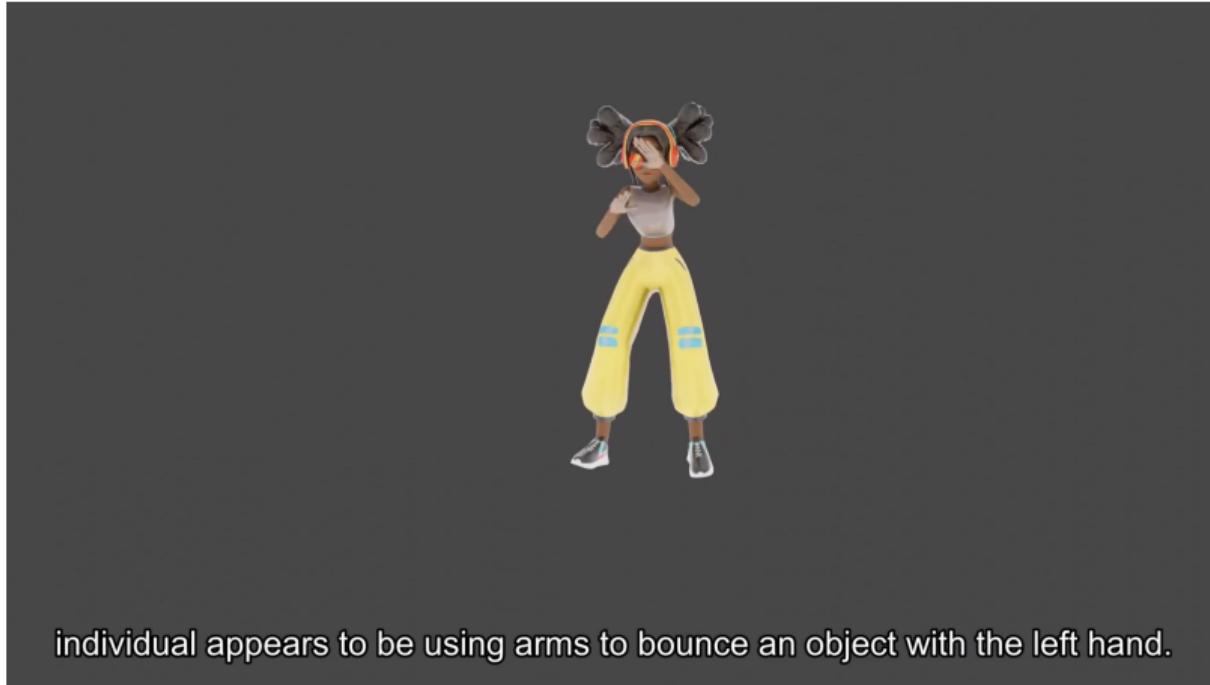


Further Application: Arbitrary Long Motion Generation



Zhang et al. *InfiniMotion: Mamba Boosts Memory in Transformer for Arbitrary Long Motion Generation*

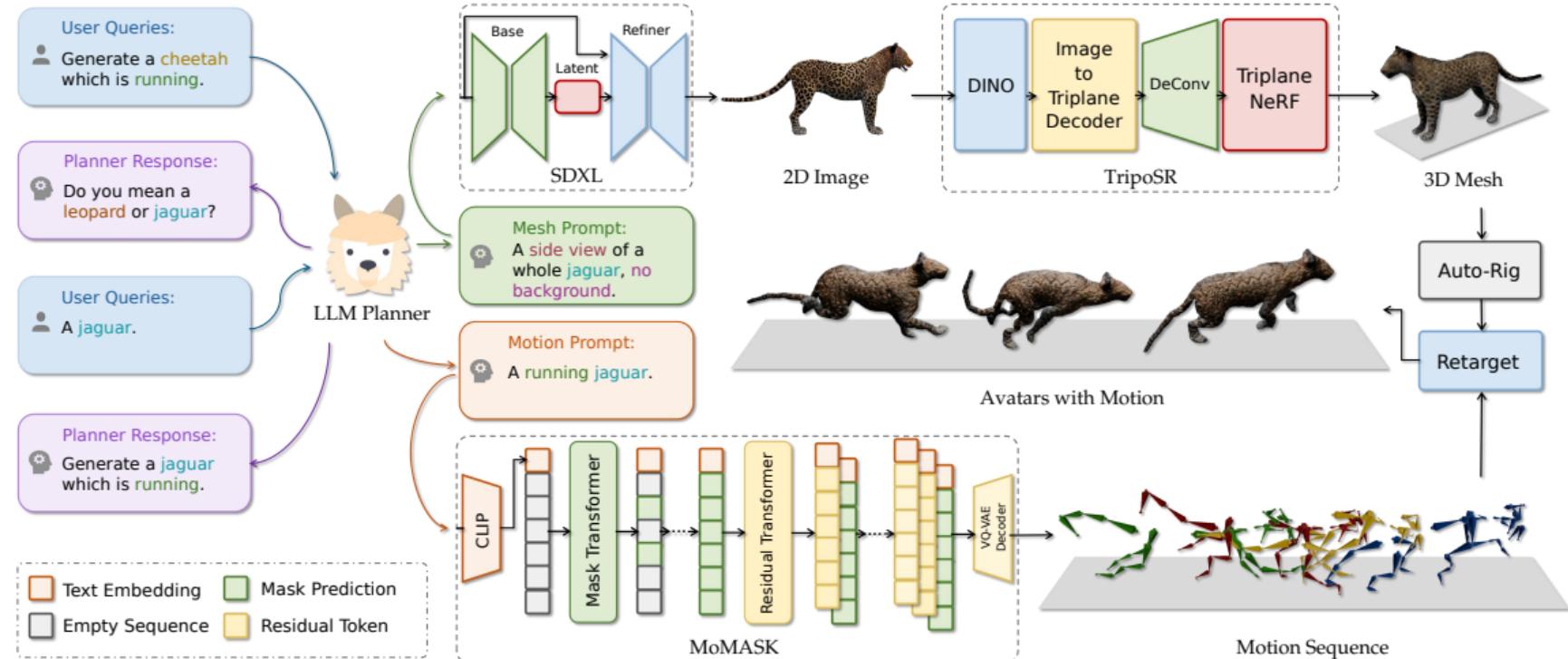
Further Application: Arbitrary Long Motion Generation



individual appears to be using arms to bounce an object with the left hand.

Zhang et al. *InfiniMotion: Mamba Boosts Memory in Transformer for Arbitrary Long Motion Generation*

Further Application: Dynamic Avatar Generation



Further Application: Dynamic Avatar Generation



A demon girl with blue hair is running fast.



A demon girl with blue hair spins rapidly.



A demon girl with blue hair holds her left foot with her left hand.



Luigi is bending his leg.



Luigi is doing kung fu.



Luigi is rolling his hands.

Further Application: HOI & HSI



Jiang et al. *TRUMANS: Scaling Up Dynamic Human-Scene Interaction Modeling* (CVPR 2024)

Limitation on Text-to-Motion Generation

- Motion generation requires much stronger conditions for controllable and precise generation, such as trajectory, video, and even joint-level control.
- Current methods still produce unrealistic joint angles and jitters. Future work should focus on generating physically compatible motions for broader real-world applications, such as robot manipulation.
- Despite the existence of works like Motion Avatar and OmniMotionGPT for animal motion generation, significant challenges remain due to the lack of datasets and the variability in animal motion representation.

Q&A

Feel free to ask any questions. Thank you!



Motion Mamba Project Website