

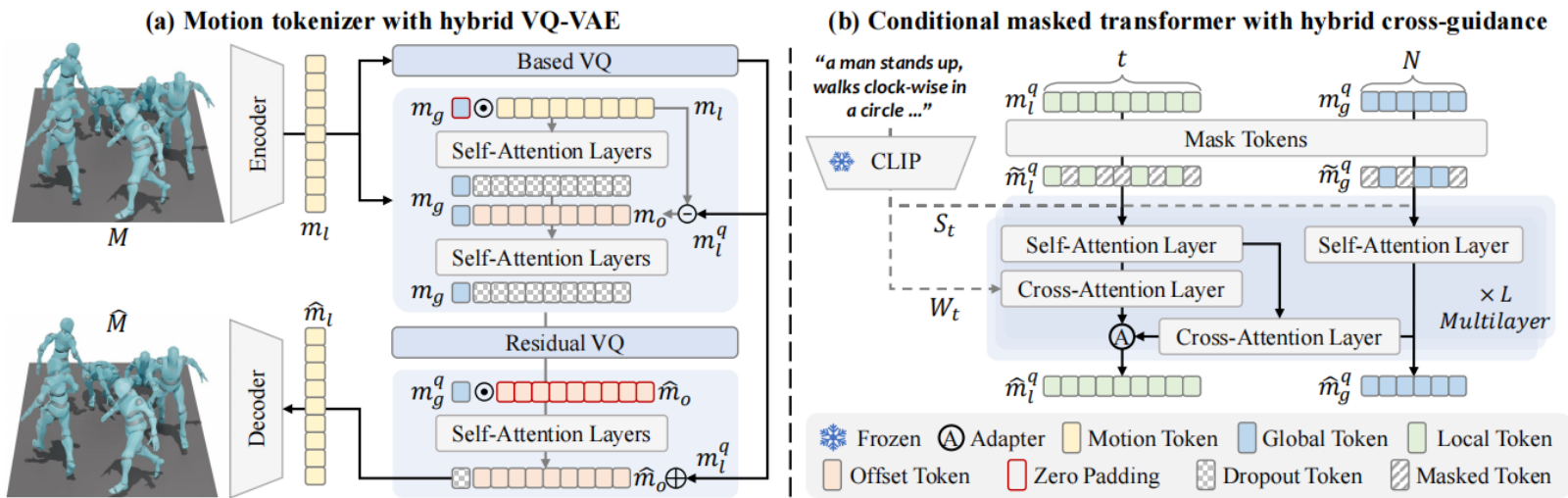
Generative Masked Text-to-Motion Model with Hybrid Vector Quantization

Jiaqi Zhang, Jiajun Wang, Fanglue Zhang, Miao Wang

Frontiers of Computer Science, DOI: [10.1007/s11704-025-50904-0](https://doi.org/10.1007/s11704-025-50904-0)

Problems & Ideas

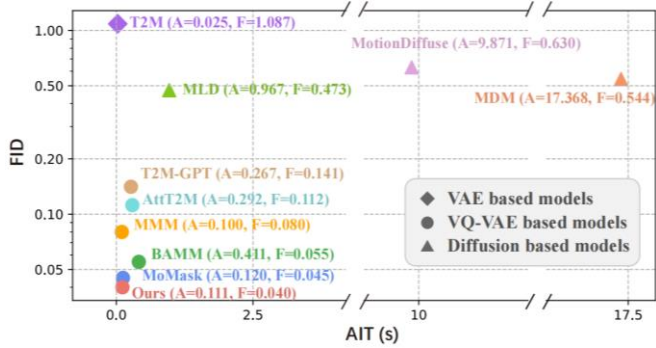
- Problems of conventional Text-Driven Motion Generation:
 - Methods based on diffusion models exhibit faster generation speeds, yet tend to compromise on generating quality.
 - Methods based VQ-VAE achieve higher generation quality but necessitate a substantial increase in the length of token sequences.
- Ideas: A hybrid VQ-VAE utilizes RVQ for global features to mitigate errors from the basic VQ-VAE, ensuring high-quality generation without substantially increasing token length.



Left Side: Displays the hybrid VQ-VAE, which decomposes motion into global and local components. The local motion is quantized using a based VQ to preserve details, while the global motion is reconstructed through a RVQ to compensate for the inaccuracies introduced by the VQ of the local components. Right Side: Featuring a conditional masked transformer guided by text and a global motion cross-attention module.

Main Contributions

- Contributions:
 - An advanced motion quantization approach (Hybrid VQ-VAE) that separates motion into global and local components, with Residual VQ applied to global motion and basic VQ-VAE to local motion;
 - The Hybrid VQ-VAE improves the quality of motion quantization while maintaining a minimal increase in code length;
 - A conditional masked transformer model that leverages high-quality reconstructions from a hybrid VQ-VAE and employs mixed cross-guidance to accurately predict global and local motion token labels.



Methods	FID ↓	R-Precision ↑			MM-Dist ↓	Diversity →	MModality ↑
		Top-1 ↑	Top-2 ↑	Top-3 ↑			
Real	0.002±.000	0.511±.003	0.703±.003	0.797±.002	2.974±.008	9.503±.065	-
TM2T [35]	1.501±.017	0.424±.003	0.618±.003	0.729±.002	3.467±.011	8.589±.076	2.424±.093
T2M [5]	1.087±.021	0.455±.003	0.636±.003	0.736±.002	3.347±.008	9.175±.083	2.219±.074
MDM [6]	0.544±.044	0.320±.005	0.498±.004	0.611±.007	5.566±.027	9.559±.086	2.799±.072
MotionDiffuse [8]	0.630±.001	0.491±.001	0.681±.001	0.782±.001	3.113±.001	9.410±.049	1.553±.042
MLD [13]	0.473±.013	0.481±.003	0.673±.003	0.772±.002	3.196±.010	9.724±.082	2.413±.079
Fg-T2M [7]	0.243±.019	0.492±.002	0.683±.003	0.783±.002	3.109±.007	9.278±.072	1.614±.049
M2DM [34]	0.352±.005	0.497±.003	0.682±.002	0.763±.003	3.134±.010	9.926±.073	3.587±.072
T2M-GPT [14]	0.116±.004	0.491±.003	0.680±.003	0.775±.002	3.118±.011	9.761±.081	1.856±.011
AttT2M [33]	0.112±.006	0.499±.003	0.690±.002	0.786±.002	3.038±.007	9.700±.090	2.452±.051
MotionLCM [9]	0.304±.012	0.502±.003	0.698±.002	0.798±.002	3.012±.007	9.607±.066	2.259±.092
MMM [15]	0.089±.005	0.515±.002	0.708±.002	0.804±.002	2.926±.007	9.577±.050	1.226±.035
MoMask [16]	0.045±.002	0.521±.002	0.713±.002	0.807±.002	2.958±.008	9.636±.068	1.241±.040
BAMM [19]	0.055±.002	0.522±.003	0.715±.003	0.808±.003	2.936±.077	9.636±.009	1.732±.055
HyT2M	0.040±.002	0.556±.003	0.749±.002	0.839±.002	2.723±.008	9.607±.057	1.042±.038

The left side shows a comparison of generation quality and inference speed of various methods on the HumanML3D test set, while the right side presents multiple quantitative metrics for these methods on the same set.