

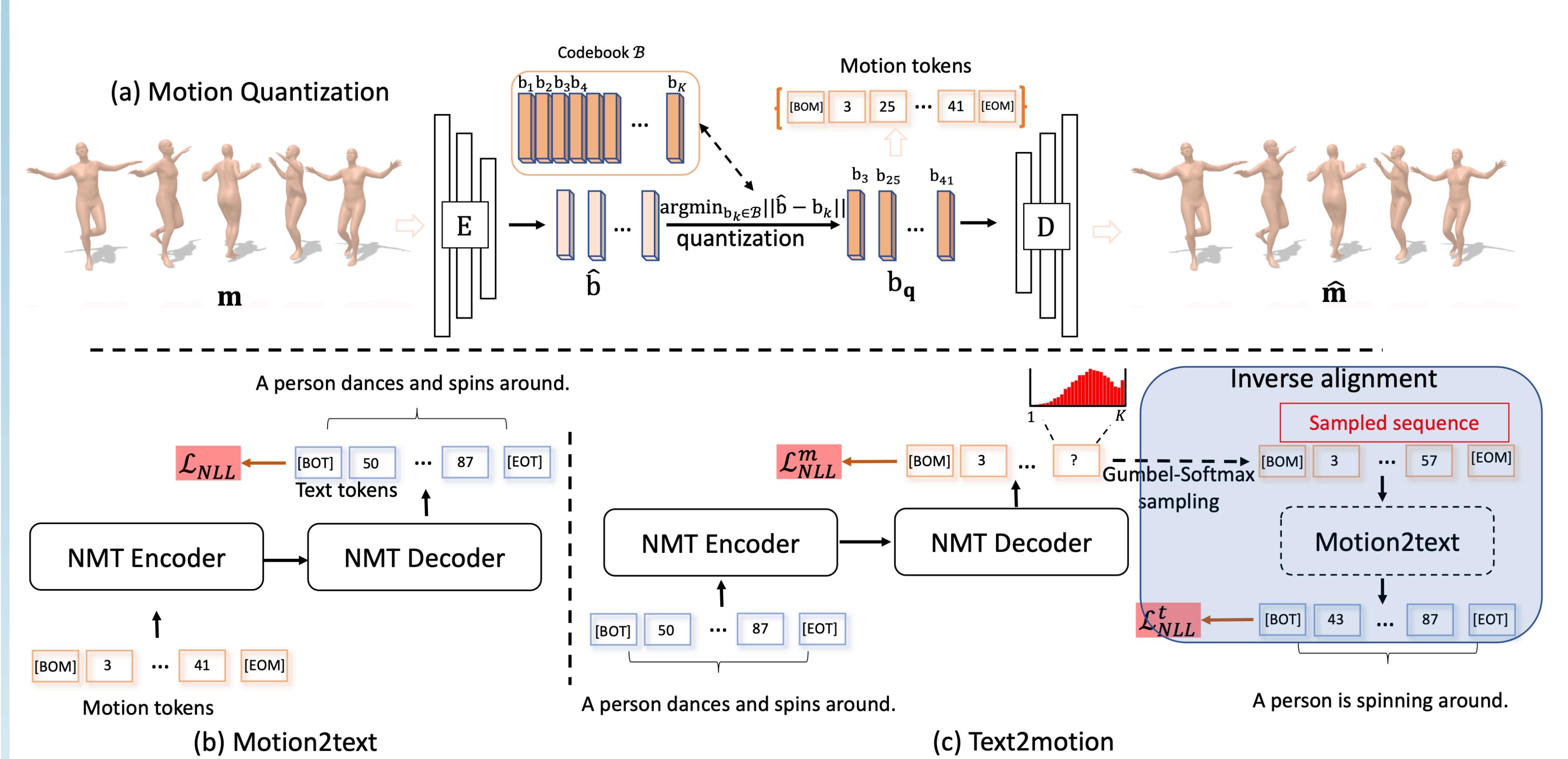
Introduction

Goal: Learning the mutual mappings between texts and motions reciprocally.

Contributions:

- We propose a novel discrete and compact representation of human motion, motion token. This provides one level playing ground when considering both motions and texts.
- Neural machine translator is adopted for stochastic text2motion and motion2text modeling.
- Motion2text is integrated into the training process of text2motion for inverse alignment.
- Extensive experiments and state-of-the-art performance on both two tasks.

Method



$$\text{Training Quantization: } \mathcal{L}_{vq} = \|\hat{\mathbf{m}} - \mathbf{m}\|_1 + \|\text{sg}[E(\mathbf{m})] - \mathbf{b}_q\|_2^2 + \beta \|\mathbf{E}(\mathbf{m}) - \text{sg}[\mathbf{b}_q]\|_2^2$$

After quantization, we can transform each motion sequence into a discrete token sequence $s \in \{1, \dots, |\mathcal{B}|\}^T$, accompanying with its text token sequence $x \in \{1, \dots, |\mathcal{V}|\}^N$, where \mathcal{V} is the word vocabulary, and the lengths of motion and text sequence are T and N respectively.

The models of motion2text θ , and text2motion ϕ are optimized as follow:

$$\text{Training Motion2text: } \mathcal{L}_{m2t} = - \sum_{i=0}^{N-1} \log p_\theta(x_i | x_{<i}, s)$$

$$\text{Training Text2motion: } \mathcal{L}_{t2m} = - \left(\sum_{i=0}^{T-1} \log p_\phi(s_i | s_{<i}, x) + \sum_{i=0}^{N-1} \log p_\theta(x_i | x_{<i}, \hat{s}) \right)$$

Dataset

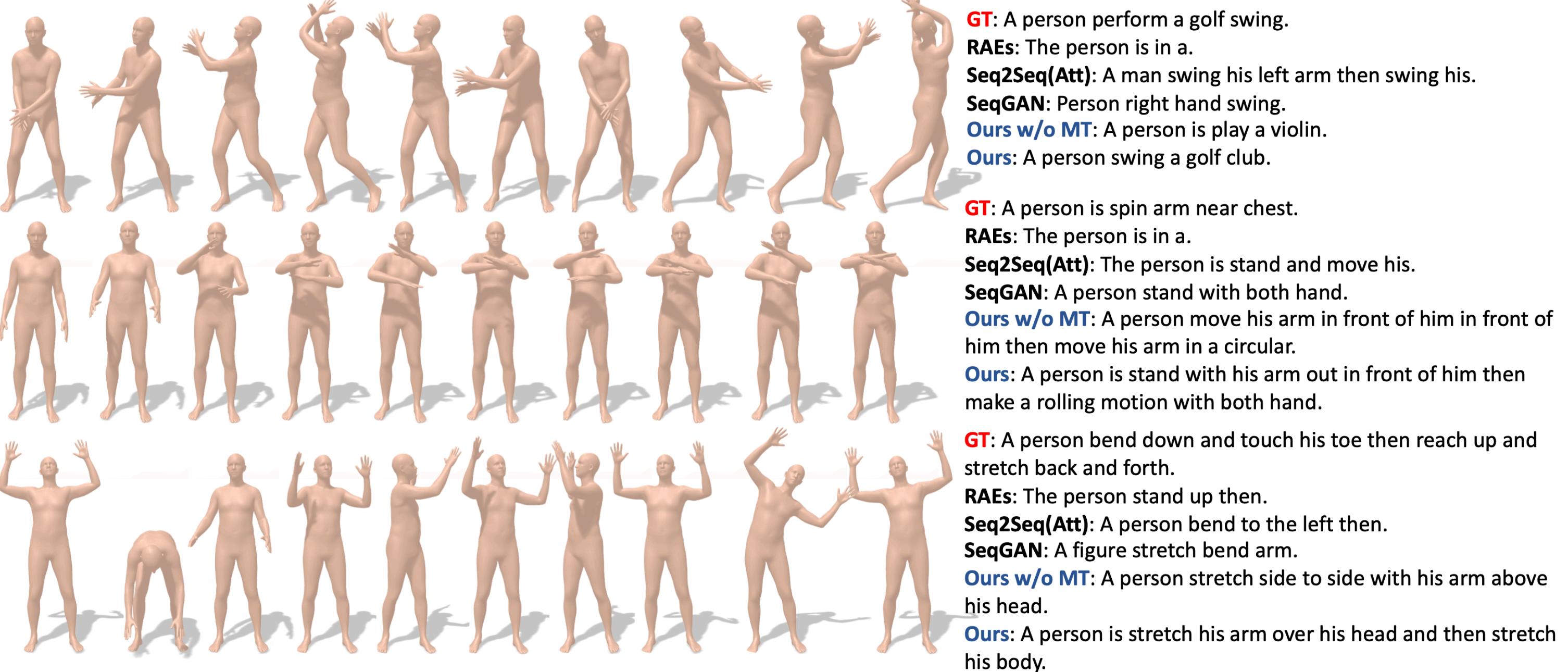
HumanML3D[1] is a large 3D human motion dataset that covers a broad range of human actions such as locomotion, sports, and dancing. It consists of **14,616** motions and **44,970** text descriptions. Each motion clip comes with at least 3 descriptions.

KIT-ML[2] contains **3,911** 3D human motion clips and **6,278** text descriptions. For each motion, the corresponding number of text descriptions ranges from one to four.

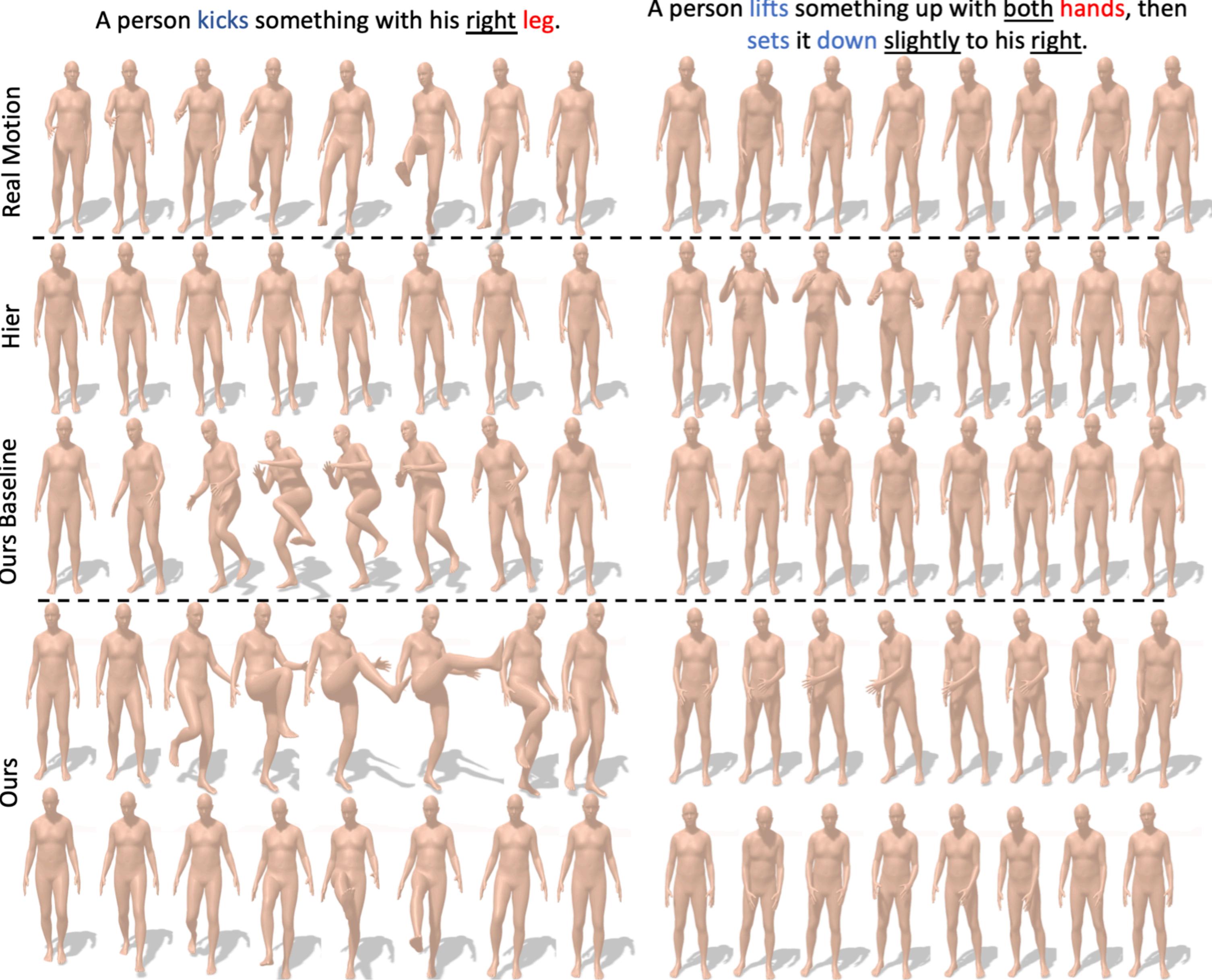
Both datasets are split into training, testing and validation sets with ratio of 0.8:0.15:0.05.

Experiments & Results

Qualitative Results of Motion2text:



Qualitative Results of Text2motion:



Quantitative Results of Motion2text:

Datasets	Methods	R Precision↑			MM Dist↓	Bleu@1↑	Bleu@4↑	Rouge↑	Cider↑	BertScore↑
		Top 1	Top 2	Top 3	-	-	-	-	-	-
Human	Real Desc	0.523	0.725	0.828	2.901	-	-	-	-	-
	RAEs[3]	0.100	0.184	0.261	6.337	33.3	10.2	37.5	22.1	10.7
	Seq2Seq(Att)	0.436	0.611	0.706	3.447	51.8	17.9	46.4	58.4	29.1
	SeqGAN[4]	0.332	0.457	0.532	4.895	47.8	13.5	39.2	50.2	23.4
ML3D	Ours w/o MT	0.483	0.678	0.783	3.124	59.5	21.2	47.8	68.3	34.9
	Ours	0.516	0.720	0.823	2.935	61.7	22.3	49.2	72.5	37.8
KIT-ML	Real Desc	0.399	0.618	0.793	2.772	-	-	-	-	-
	RAEs[3]	0.034	0.063	0.106	9.364	30.6	0.10	25.7	8.00	0.40
	Seq2Seq(Att)	0.293	0.450	0.555	4.455	34.3	9.30	36.3	37.3	5.30
	SeqGAN[4]	0.109	0.345	0.425	6.283	3.12	5.20	32.4	29.5	2.20
ML3D	Ours w/o MT	0.284	0.466	0.595	3.979	42.8	14.7	39.9	60.1	18.9
	Ours	0.359	0.561	0.668	3.298	46.7	18.4	44.2	79.5	23.0

Quantitative Results of Text2motion:

Datasets	Methods	R Precision↑			FID↓	MM Dist↓	Diversity→ MModality↑
		Top 1	Top 2	Top 3	-	-	-
Human	Real motions	0.511 ^{±.003}	0.709 ^{±.002}	0.791 ^{±.002}	0.002 ^{±.000}	2.974 ^{±.007}	9.503 ^{±.065}
	Seq2Seq[5]	0.180 ^{±.002}	0.300 ^{±.002}	0.396 ^{±.002}	11.75 ^{±.035}	5.529 ^{±.061}	6.223 ^{±.058}
	Language2Pose[6]	0.246 ^{±.002}	0.387 ^{±.002}	0.486 ^{±.002}	11.02 ^{±.046}	5.296 ^{±.068}	7.676 ^{±.058}
	Text2Gesture[7]	0.165 ^{±.001}	0.267 ^{±.002}	0.345 ^{±.002}	5.012 ^{±.030}	6.030 ^{±.071}	6.409 ^{±.071}
ML3D	Hier[8]	0.301 ^{±.002}	0.422 ^{±.002}	0.555 ^{±.002}	6.532 ^{±.024}	5.012 ^{±.018}	8.332 ^{±.042}
	MoCoGAN[9]	0.037 ^{±.000}	0.072 ^{±.000}	0.106 ^{±.000}	94.41 ^{±.001}	9.643 ^{±.000}	0.462 ^{±.008}
	Dance2Music[10]	0.033 ^{±.000}	0.065 ^{±.001}	0.097 ^{±.001}	66.98 ^{±.016}	8.116 ^{±.000}	0.725 ^{±.011}
	Ours baseline(T)	0.351 ^{±.003}	0.521 ^{±.003}	0.627 ^{±.003}	1.669 ^{±.025}	4.046 ^{±.018}	9.632 ^{±.072}
KIT-ML	Ours baseline	0.351 ^{±.002}	0.526 ^{±.002}	0.633 ^{±.002}	1.739 ^{±.022}	3.965 ^{±.010}	8.651 ^{±.083}
	Ours	0.424^{±.003}	0.618^{±.003}	0.729^{±.002}	1.501^{±.017}	3.467^{±.011}	8.589^{±.076}
ML3D	Real motions	0.424 ^{±.005}	0.649 ^{±.006}	0.779 ^{±.006}	0.031 ^{±.004}	2.788 ^{±.012}	11.08 ^{±.097}
	Seq2Seq[5]	0.103 ^{±.003}	0.178 ^{±.003}	0.241 ^{±.003}	24.86 ^{±.348}	7.960 ^{±.106}	6.744 ^{±.106}
	Language2Pose[6]	0.221 ^{±.005}	0.373 ^{±.005}	0.483 ^{±.005}	6.545 ^{±.072}	5.147 ^{±.088}	9.073 ^{±.100}
	Text2Gesture[7]	0.156 ^{±.004}	0.257 ^{±.004}	0.334 ^{±.005}	12.19 ^{±.183}	6.964 ^{±.079}	9.334 ^{±.079}
ML3D	Hier[8]	0.255 ^{±.006}	0.432 ^{±.007}	0.531 ^{±.007}	5.203 ^{±.107}	4.986 ^{±.027}	5.963 ^{±.072}
	MoCoGAN[9]	0.022 ^{±.002}	0.042 ^{±.003}	0.063 ^{±.003}	82.69 ^{±.242}	10.47 ^{±.012}	3.091 ^{±.043}
	Dance2Music[10]	0.031 ^{±.002}	0.058 ^{±.002}	0.079 ^{±.002}	115.4 ^{±.240}	10.40 ^{±.010}	0.241 ^{±.004}
	Ours baseline(T)	0.260 ^{±.005}	0.426 ^{±.007}	0.538 ^{±.008}	4.628 ^{±.126}	4.835 ^{±.076}	12.16 ^{±.120}
KIT-ML	Ours baseline	0.251 ^{±.007}	0.418 ^{±.008}	0.535 ^{±.007}	4.814 ^{±.145}	4.682 ^{±.048}	10.13 ^{±.117}
	Ours	0.280^{±.005}	0.463^{±.006}	0.587^{±.005}	3.599<		