## (a) Motion Tokenizer (b) Masked Self-attention Transformer Key Key Ε 3 6 2 3 4 5 6 E Decoder **Masked Self-attention** 1 2 3 4 5 6 Query Query Codebook **Masked Self-attention** 1 2 ... 512 Encoder Bidirectional Causal Mask $M^{bc}$ Unidirectional Causal Mask $M^{uc}$ **Encoder** Conditioned / target motion token Attention not allowed $(-\infty)$ Text token CLIP Motion Token in Masked Area E End token Attention allowed (0) "A person walks in a circle"