

→ : Forward Pass

↔ : Stop Gradient

⊕ : Level-wise concatenation

Weight Shared

C^m : Motion-level Text embedding

C^a : Action-level Text embedding

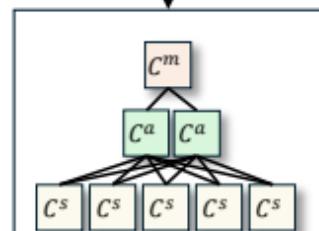
C^s : Specifics-level Text embedding

Padding

Input Text

a person takes a huge leap forward.

Graph Reasoning



Input Motion Sequence X

Residual VQ-VAE Encoder \mathcal{E}

Original Motion Tokens Y^0

