

Legend: ■ Hidden States (h) ■ [listen] Token ■ Text Token ⊕ Token + h → Decoder | ■ Visual Embeddings (V) ■ Audio Embeddings (A)

