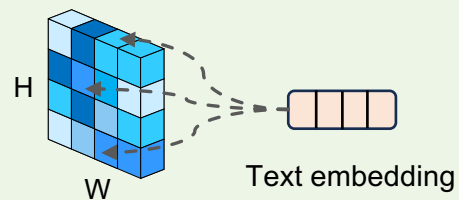
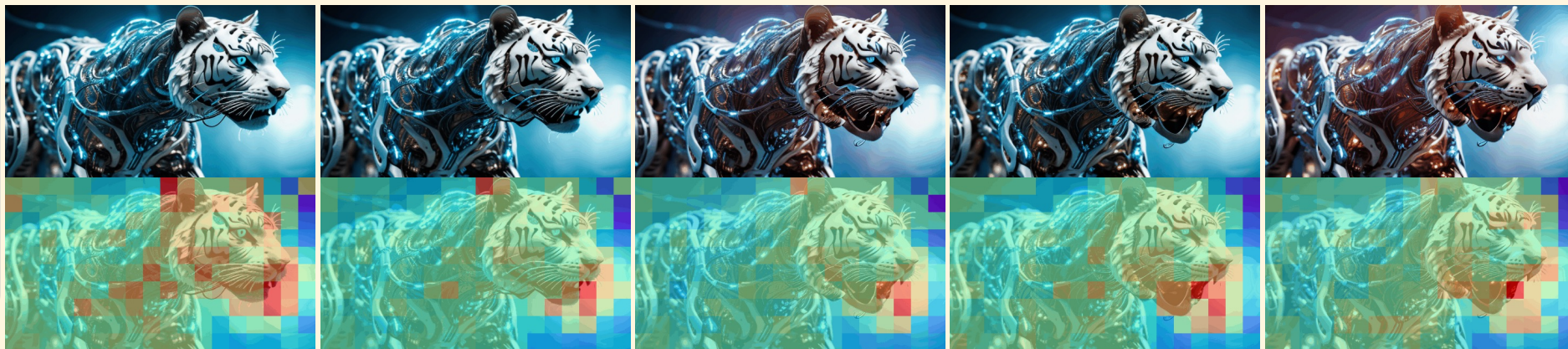
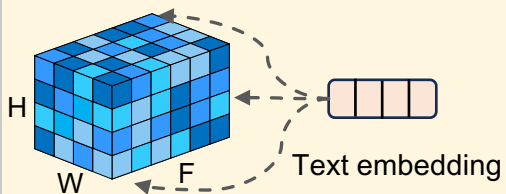


(a) Spatial 2D  
Cross-attention



A mechanical white tiger is roaring

(b) Motion-aware  
3D Cross-attention



A mechanical white tiger is roaring