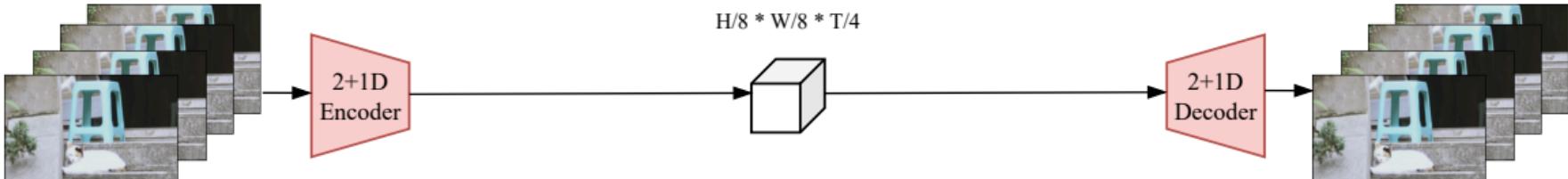
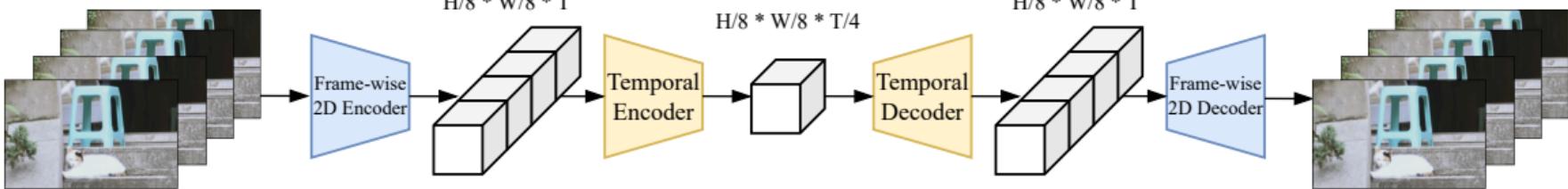
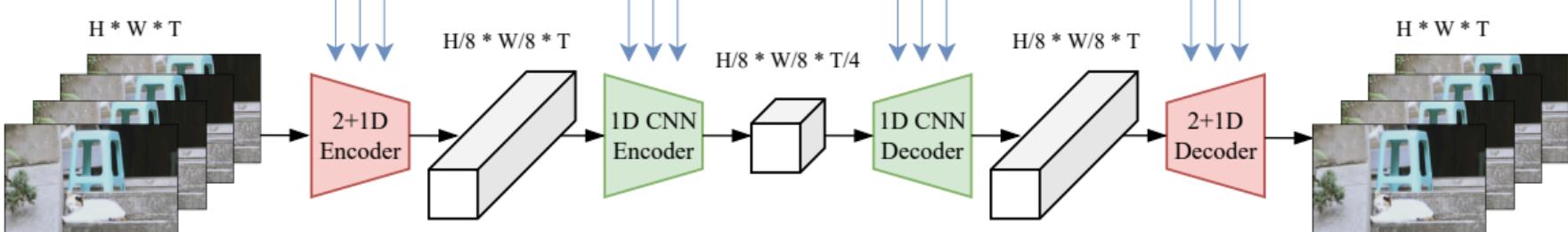


$H * W * T$ $H * W * T$ **Simultaneously modeling** $H * W * T$ $H/8 * W/8 * T$ $H/8 * W/8 * T/4$ $H/8 * W/8 * T$ $H * W * T$ **Sequential modeling**

Description: "A peaceful cat is lying on the stairs, nestled against the warm steps, basking in the soft afternoon sunlight"

 $H * W * T$ $H/8 * W/8 * T$ $H/8 * W/8 * T/4$ $H/8 * W/8 * T$ $H * W * T$ **Ours (Combine both advantages and use cross-modal modeling)**