

[N=1, T=1, H, W, C]



A kid doing a kickflip.



to my



Image-Text Pairs dataset

Video-Text Pairs dataset

[N=1, T>1, H, W, C]

Multi-Modal Massive Web (M3W) dataset [N>1, T=1, H, W, C]