

Verb-object query:
“Slicing vegetables”

Verb-object feature



Video frame I_t

Human
region
features f_t^h

Object region features f_t^o

Contextual
frame feature x_t

Object region features
from other video frames \mathcal{F}

**Region
attention**

Attended
human
feature Φ_t^h

Attended
object
feature Φ_t^o

Weakly supervised contrastive loss