



ROMA's Streaming Understanding Ability

Proactive

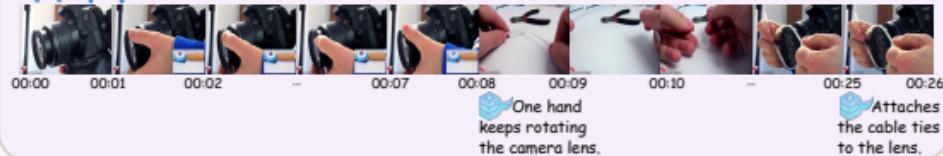
Event-Driven Alert

(Notify me when a bird pops out from behind a tree.)



Real-Time Narration

(What is occurring with the camera lens?)



Reactive

Reactive QA

(What is the woman doing?) (What is the last step to complete the soup?)



She's **cooking** ... a pile of **<eat>** tomatoes, then slicing them up.

Garnish the soup with chopped basil.