



So far, we only use *Natural Language Query*  $q_t$ :

*Sentence:* Woman holds her a lobster coffee mug.

What comes next: *Multimodal Query*  $q_m$ :


(Reference image  + Refinement text  $\overset{\text{TXT}}{\nearrow}$  )

User scribble



In the kitchen

Similar scenes



A lobster mug  
Not a real lobster

