**a.** text concept extraction

"A dog and a cow on grass."

concept identification

"A dog and a cow on grass."

text encoder

**b.** cross-modality consistency

**c.** vision concept extraction

"dog" "cow" "grass"

text-guided pooling

linear classifier

vision encoder

text representation
vision representation