

(a) Demonstration Setup



Vision-based Demo

<image> 0%



<image> 25%



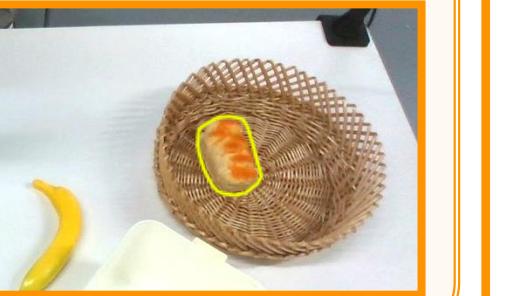
<image> 50%



<image> 75%



<image> 100%



Text-based Demo

<text> 25%
"move towards
the bread"

<text> 50%
"grab the bread"

<text> 75%
"move the bread
towards the basket"

<text> 100%
"move away from
the basket"



(c) Answerability Augmentation

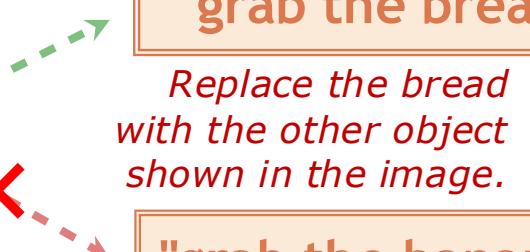
Text-based Demo

"grab the bread"

Replace the bread
with the other object
shown in the image.

"grab the banana"

o^i



o^i

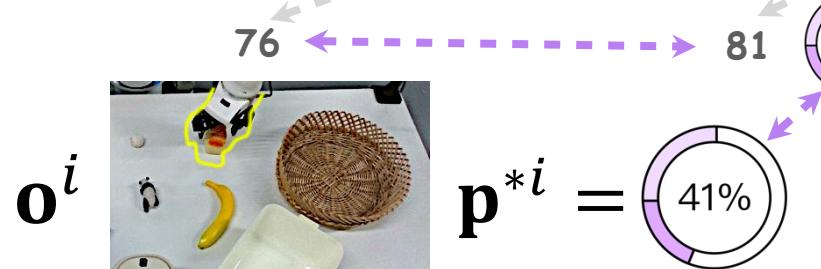


Vision-based Demo

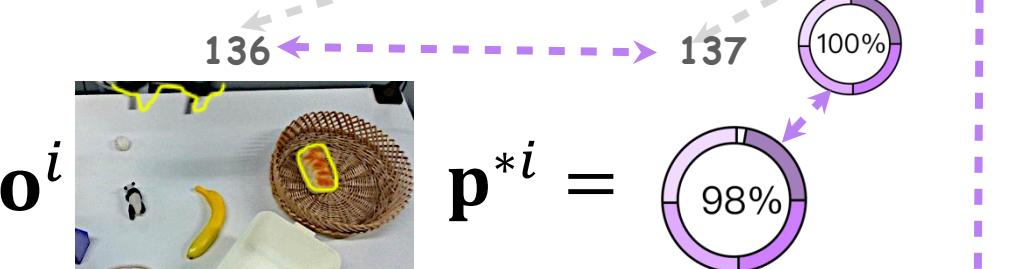
Replace the bread
to the apple



(b) Observation Sampling

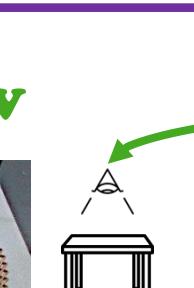


interval sampling



boundary sampling

same-view



Vision-based Demo

cross-view



demonstration-observation correspondence