

## Textual Input $X$

The color of the banana is ...

Original Visual  
Input  $V$



$\text{logits}(y|x, v)$

|        |              |
|--------|--------------|
| Black  | <b>16.72</b> |
| Dark   | 14.45        |
| Yellow | 11.27        |
| Green  | 12.30        |

Distorted Visual  
Input  $V'$



$\text{logits}(y|x, v')$

|        |                |
|--------|----------------|
| Black  | 12.57          |
| Dark   | 11.84          |
| Yellow | <b>14.74</b> ↑ |
| Green  | 13.57 ↑        |

## Distorted Visual Input

