

*High similarity*

Original image  $v$



$\succ =$

Modified image  $\hat{v}_1$



$\succ =$

*Low similarity*

Modified image  $\hat{v}_2$



Input prompt

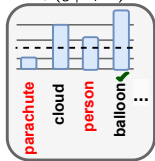
"In the sky there is a ..."



"... **balloon**"

Unlikely tokens generated from ordinarily modified images are penalized.

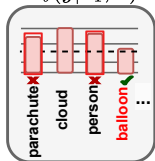
$P_\theta(y|v, x)$



(A.)

$\succ =$

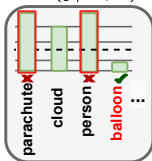
$P_\theta(y|\hat{v}_1, x)$



(B.)

$\succ =$

$P_\theta(y|\hat{v}_2, x)$



(C.)