

Image input: \mathbf{I}



Text input: \mathbf{T}

men in suits are
walking in a park

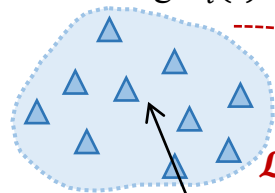
Image
Encoder

$E_i(\cdot)$

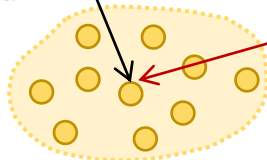
Text
Encoder

$E_t(\cdot)$

Original image
embedding $E_i(\mathbf{I})$



Original text
embedding $E_t(\mathbf{T})$

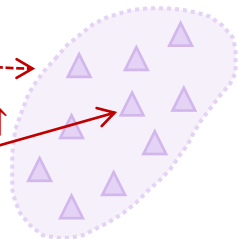


Attack on visual modality

$\mathcal{L}(E_i(\mathbf{I}'), E_i(\mathbf{I})) \uparrow$

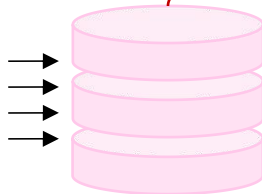
$\mathcal{L}(E_i(\mathbf{I}'), E_t(\mathbf{T})) \uparrow$

Adversarial image
embedding $E_i(\mathbf{I}')$



Attack on textual modality

BERT-Attack-
based generator



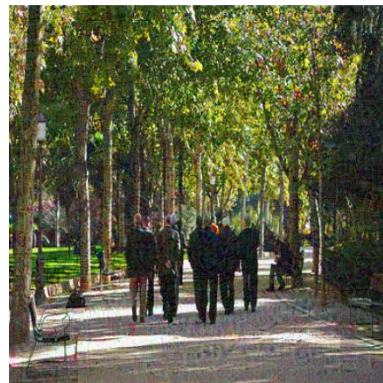
Adversarial
text library

Select sample

$\mathcal{L}(E_t(\mathbf{T}'), E_i(\mathbf{I}')) \uparrow$

$\mathcal{L}(E_t(\mathbf{T}'), E_t(\mathbf{T})) \downarrow$

Adversarial image: \mathbf{I}'



Adversarial text: \mathbf{T}'

men in suits are
walking in a beautiful



Downstream tasks



...