

Input

 x w_1
 w_2
 w_3
 \vdots
 w_{i-1}
 w_i
 w_{i+1}
 \vdots
 w_n

Randomized substitution

Embedding
space

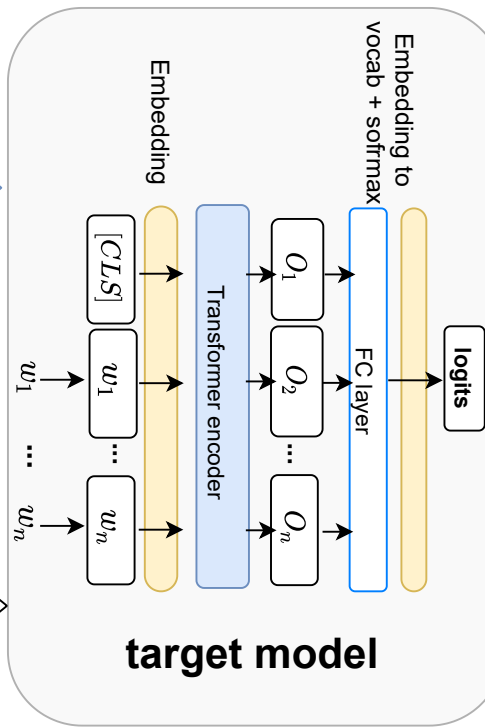
WordNet

synonym
set \hat{w}_i^1
 \hat{w}_i^2
 \hat{w}_i^3
 \vdots
 \hat{w}_i^d

converted texts

 x_1 \hat{w}_1^1
 w_2
 \hat{w}_3^2
 \vdots
 \hat{w}_{i-1}^3
 \hat{w}_i^d
 w_{i+1}
 \vdots
 \hat{w}_n^d x_2 w_1
 \hat{w}_2^1
 \hat{w}_3^2
 \vdots
 \hat{w}_{i-1}^d
 \hat{w}_i^3
 \hat{w}_{i+1}^3
 \vdots
 w_n x_3 \hat{w}_1^2
 \hat{w}_2^3
 w_3
 \vdots
 \hat{w}_{i-1}^1
 \hat{w}_i^d
 \hat{w}_{i+1}^1
 \vdots
 \hat{w}_n^3 x_k \hat{w}_1^d
 w_2
 \hat{w}_3^1
 \vdots
 \hat{w}_{i-1}^3
 \hat{w}_i^2
 w_{i+1}
 \vdots
 \hat{w}_n^2

Query



Vote & Detection

 $f(x_1)$ $f(x_2)$ $f(x_3)$ \vdots $f(x_k)$

$$\operatorname{argmax} \sum_{i=1}^k f(x_i)$$

RS&V label

adversarial
if unequal $f(x)$

logits

$$\operatorname{argmax} f(x)$$

original label