

**A Original  $x$**

It is great for kids. ( $y=\text{pos}$ ,  $f(x)=\text{neg}$ )

**B Generated  $\hat{x}$**

negation

It is **not** great for kids.  
It is great for **kids**→**no one**.

delete

It is great **for kids**.

lexical

It is **great**→**terrible** for kids.  
It is **great**→**meant** for kids.

**C Use scenarios**  
 $w/\text{ranks on } x \rightarrow \hat{x}$

**Training**

*Diversity*

**Evaluation**

*Diversity + changed label*

**Explanation**

*Abnormality*

*Targeted perturbation*