

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** \rightarrow **no one**.

delete

It is great **for kids**.

lexical

It is **great** \rightarrow **terrible** for kids.
It is **great** \rightarrow **meant** for kids.

C Downstream scenarios
 $w/\text{ selection on } x \rightarrow \hat{x}$

Training / augmentation

Diversity

Evaluation / Contrast set

Diversity + changed label

Counterfactual explanation

Abnormality

Interactive explanation

Perturbation type/subphrases