# (a) Overview

Intact input



Explicit perturbation

Shared

Shared

Implicit perturbation

$\mathbf{z}_{\text{Intact}}$

$\mathbf{z}_{\text{Biased}}$

$\mathbf{z}_{\text{Biased}}$

$y$

**Non-uniform distribution across classes**

$U$

**Uniform distribution across classes**

# (b) Debiasing Contrastive Loss

Input batch of n examples: Output: 2n + 1

**Contradiction**  **Entailment**  **Neutral**  **Perturbed**

Similar  Similar  Similar

Dissimilar

$U$

=

Dummy

Similar