

Input image  $X$



Trans.

Normalized image  $\tilde{X}$



Feat.

Visual feature  $V$



Seq.

Contextual feature  $H$



Pred.

Prediction  $Y$

UNITED