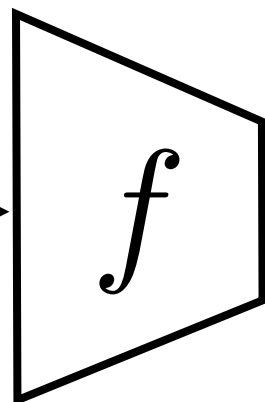


## Observation

A profile photo  
of a robin,  
facing left.

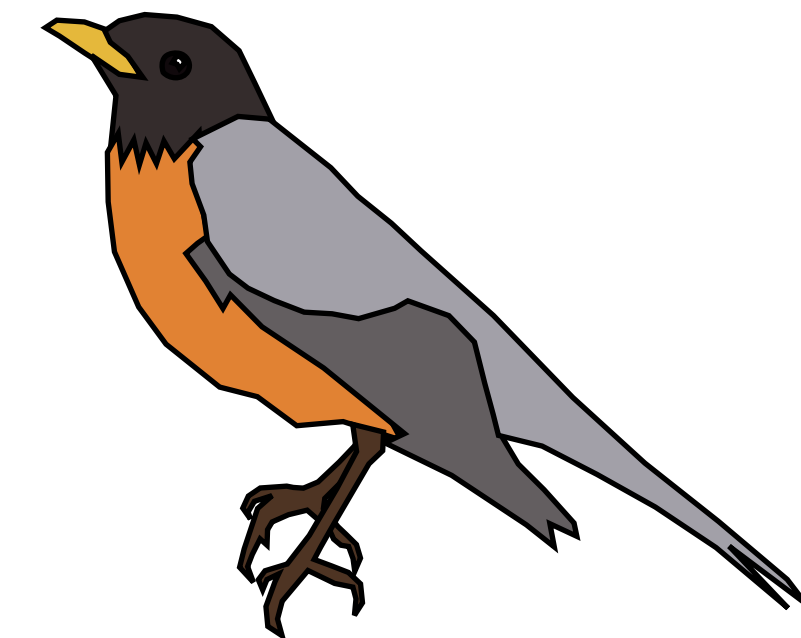
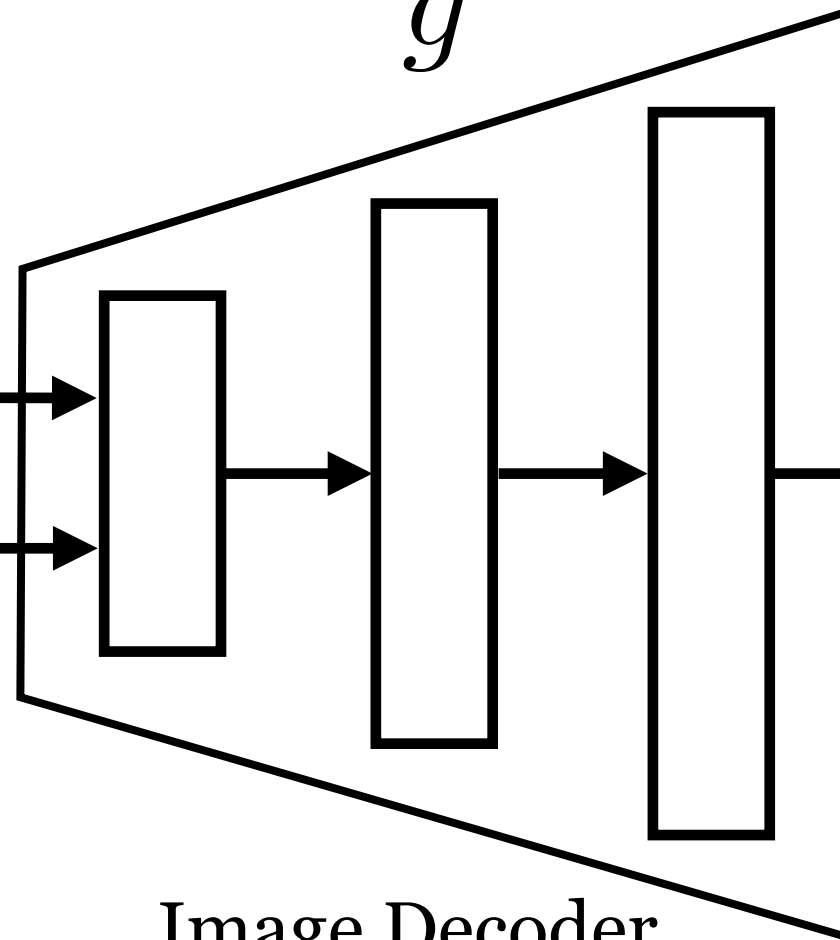
$\mathbf{t}^{(i)}$

## Text Encoder



$\mathbf{z}^{(i)} \sim \mathcal{N}(0, 1)$

$g$



Prediction

$\hat{\ell}^{(i)}$

Image Decoder