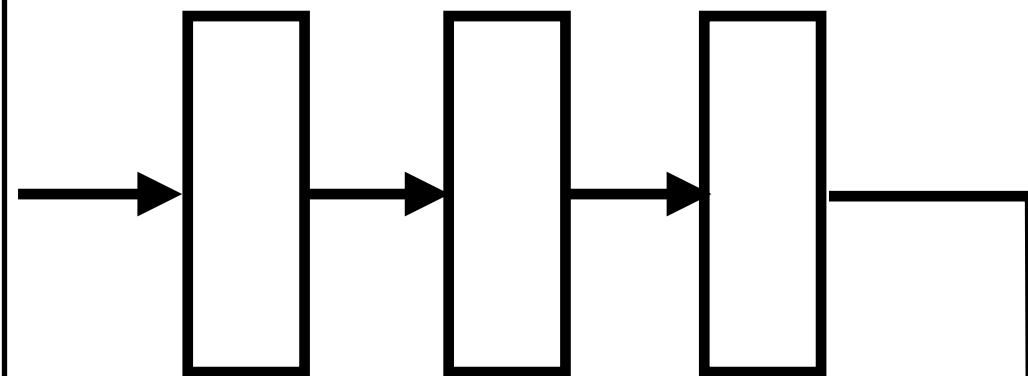


Observation

"A wide-eyed cat
on the lookout
for food"

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



Text Encoder

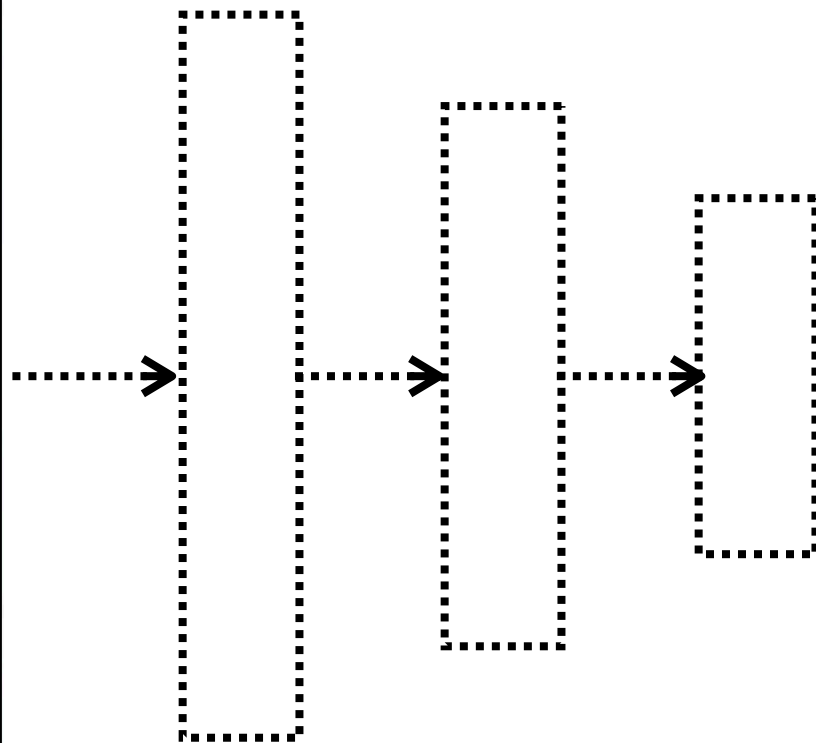


Image Encoder



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

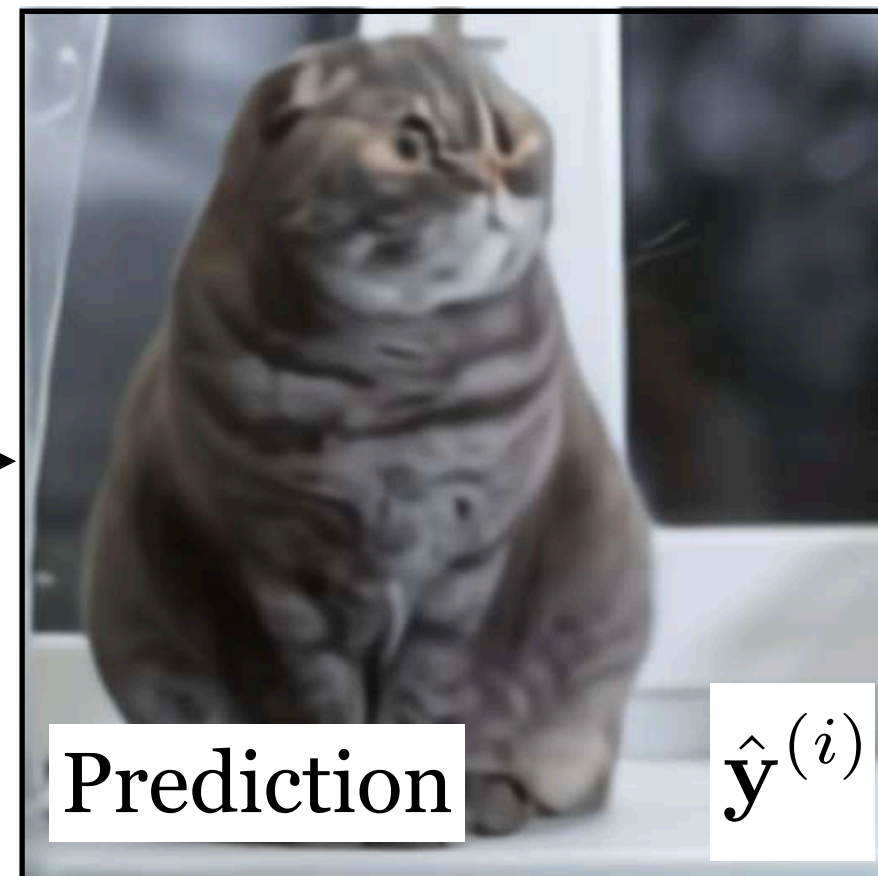
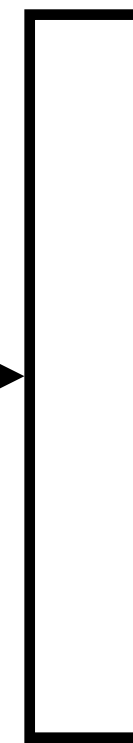
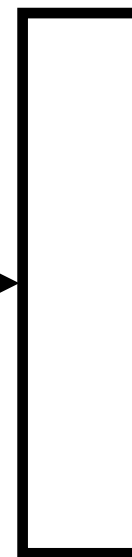
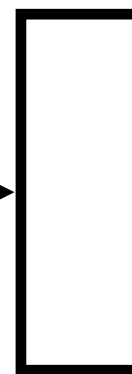


Image Decoder