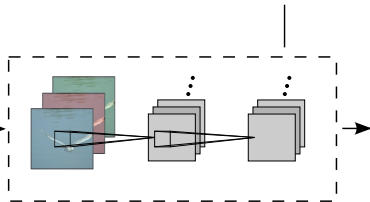


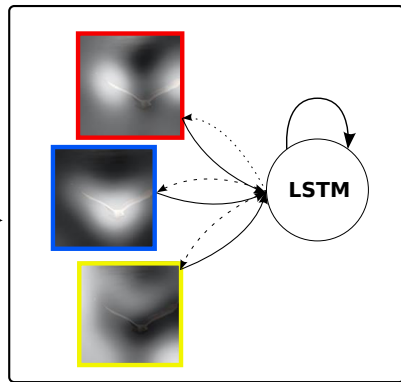


1. Input Image

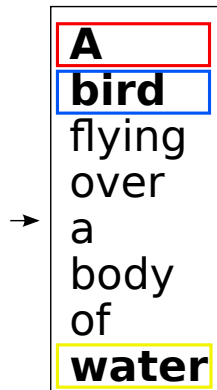
14x14 Feature Map



2. Convolutional Feature Extraction



3. RNN with attention over the image



4. Word by word generation