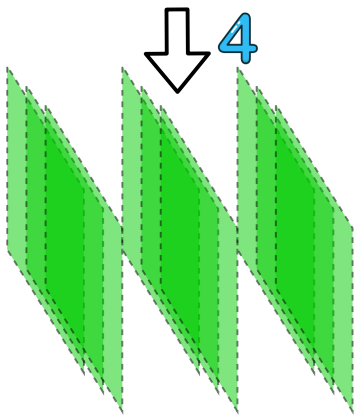


head of a cat statuette,  
a pair of standing ears,  
four dark brown thin legs, a cat sitting.

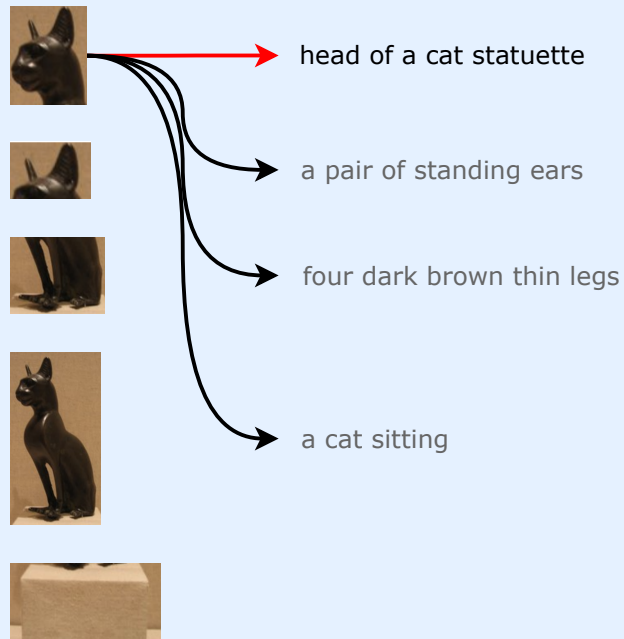


ResNet, Fully-connected layers, GRUs



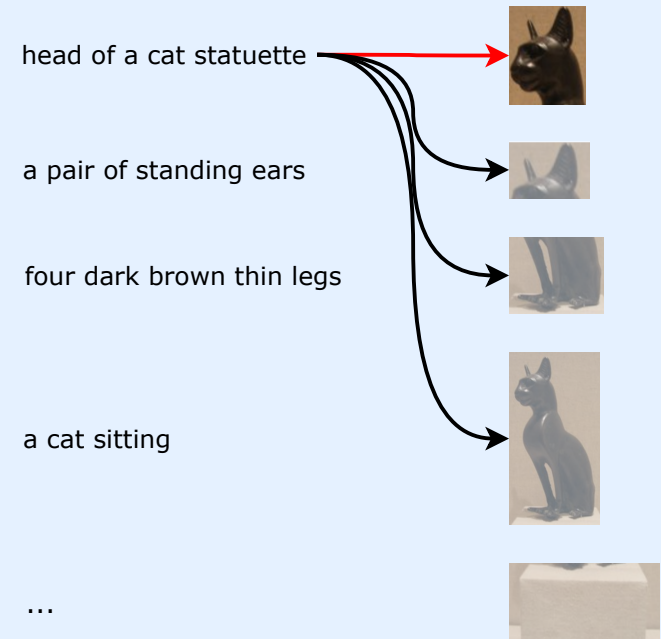
representations of  
each phrase and its  
corresponding image  
fragment now in a  
common space

attend to noun phrases with respect  
to each image feature



align targets

attend to image features with respect  
to each noun phrase



Calculate mutual attention  
between these noun  
phrases and the image  
fragments obtained from  
*faster-rcnn*

**Stacked Cross Attention**

Calculate similarity  
between image  
fragment and phrases



**Ranked recall**