

Input Image

Faster
R-CNN

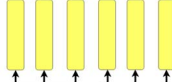
FC

Self-Attention Module



$\{r_1, r_2, \dots, r_k\}$

Transformer



Pool

img_emb: i_0



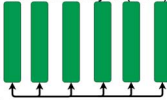
Input Sentence

Giraffes and ostriches are
sharing a grassy area.

{ '[CLS]', 'Giraffes', 'and', 'ostriches', 'are',
'sharing', 'a', 'grassy', 'area', '.', '[SEP]' }

Token Embedding +
Position Embedding +
Segment Embedding

Transformer 12 layers (BERT)



$\{e_1, e_2, \dots, e_n\}$

Self-Attention Module



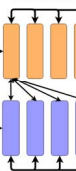
Pool

cap_emb: c_0



Cross-Attention Module

Transformer



$\{r_{c1}, \dots, r_{ck}, e_{c1}, \dots, e_{cn}\}$

Pool

img_emb: i_1



cap_emb: c_1



1d CNN
& Pool