



MSCOCO

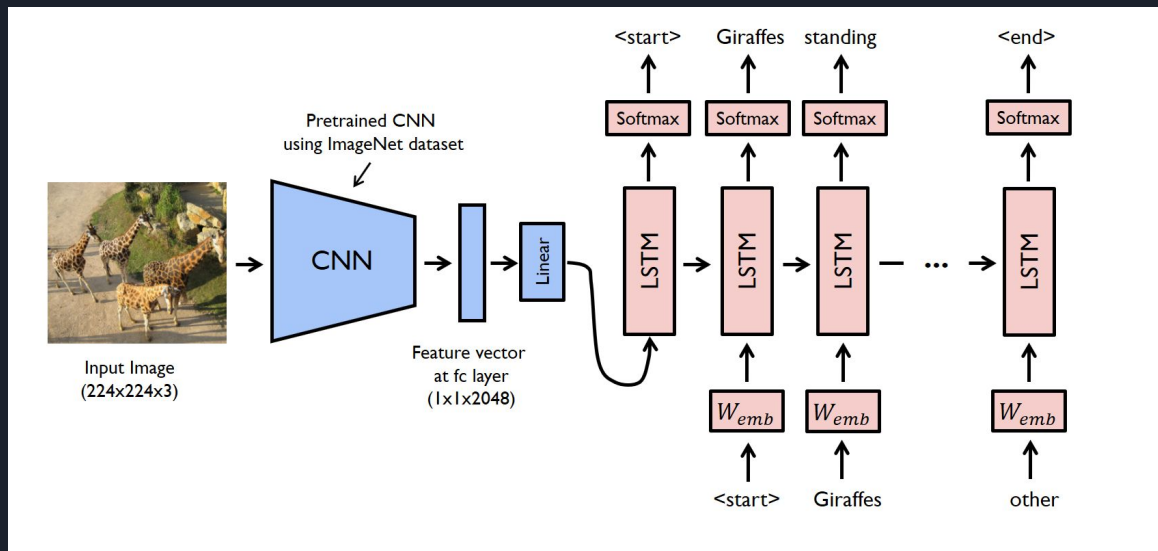
Team 8 Fortune-teller

Ping Zhu, Yang Hu, Zhichao Yang

Model — Overview

Encoder
(resnet 152)

Decoder
(LSTM)



Model — Knowing when to look

Knowing When to Look: Adaptive Attention via A Visual Sentinel for Image Captioning

Jiasen Lu^{2*†}, Caiming Xiong^{1†}, Devi Parikh³, Richard Socher¹

¹Salesforce Research, ²Virginia Tech, ³Georgia Institute of Technology
jiasenlu@vt.edu, parikh@gatech.edu, {cxiong, rsocher}@salesforce.com

Encoder

Use Resnet 152 to extract the image features V and V_g where v_g is used to form x_t .

$$A = \{a_1, \dots, a_k\}, a_i \in \mathcal{R}^{2048}$$

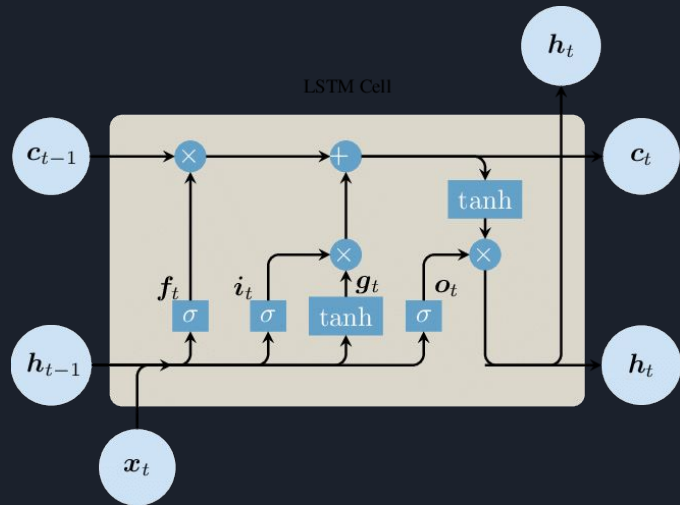
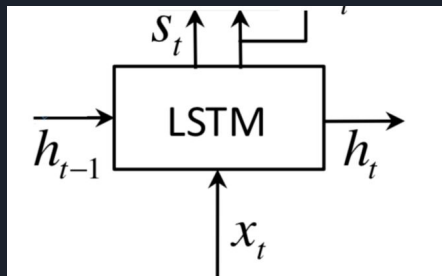
$$a^g = \frac{1}{k} \sum_{i=1}^k a_i$$

$$\begin{aligned} v_i &= \text{ReLU}(W_a a_i) \\ v^g &= \text{ReLU}(W_b a^g) \end{aligned}$$

Knowing when to look — Decoder

Modified LSTM cell to
extract sentinel

$$\begin{aligned} \mathbf{g}_t &= \sigma(\mathbf{W}_x \mathbf{x}_t + \mathbf{W}_h \mathbf{h}_{t-1}) \\ \mathbf{s}_t &= \mathbf{g}_t \odot \tanh(\mathbf{m}_t) \end{aligned}$$



Knowing when to look — Decoder

Attention

$$z_t = w_h^T \tanh(W_v V + (W_g h_t) \mathbf{1}^T)$$

$$\alpha_t = \text{softmax}(z_t)$$

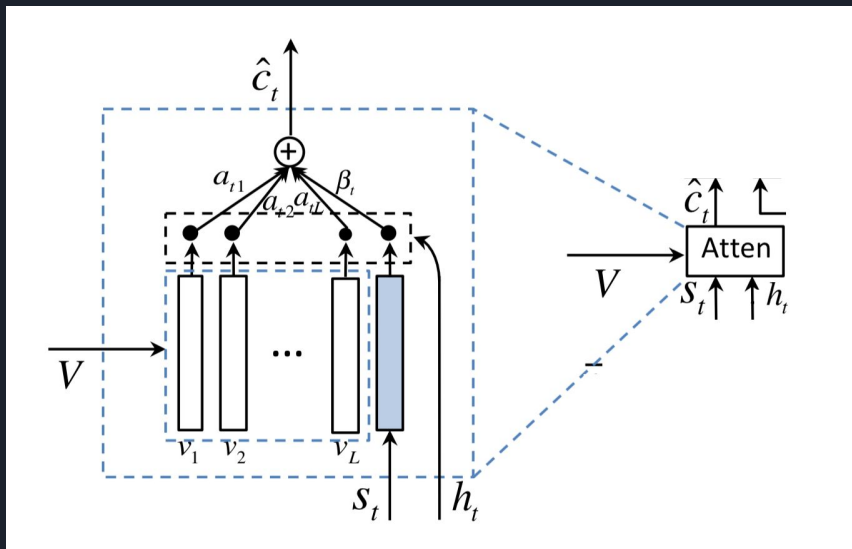
$$\hat{\alpha}_t = \text{softmax}([z_t; w_h^T \tanh(W_s s_t + (W_g h_t))])$$

$$c_t = \sum_{i=1}^k \alpha_{ti} v_{ti}$$

$$\hat{\alpha}_t \in \mathcal{R}^{k+1}$$

$$\beta_t = \alpha_t[k+1]$$

$$\hat{c}_t = \beta_t s_t + (1 - \beta_t) c_t$$

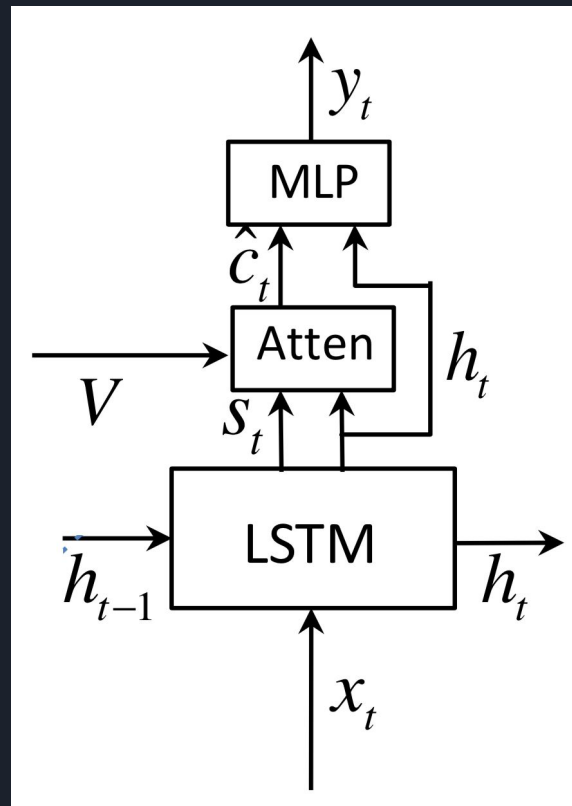


Knowing when to look — Decoder

Predict & Overview

$$p_t = \text{softmax}(\mathbf{W}_p(\hat{\mathbf{c}}_t + \mathbf{h}_t))$$

batch_size x # x len(vocab)





Performance

	CIDEr	Bleu 1	Bleu 2	Bleu 3	Bleu 4	ROUGE_L
C5	0.67169	0.66128	0.4769	0.3254	0.21609	0.48193
C40	0.7045	0.8432	0.7246	0.5873	0.4528	0.61058

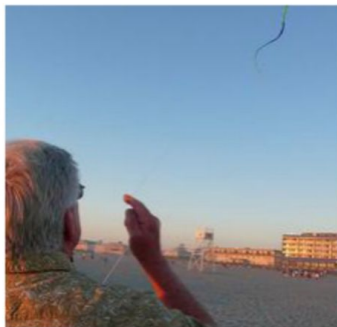
Good Predictions

GROUND TRUTH CAPTION:

A man flying a kite on the beach.
An older man watches a kite fly from across a body of water.
An older gentleman flies a kite on the beach.
A man that is standing in the sand flying a kite.
An old man flying a kite high up in the air.

PREDICTED CAPTION:

a man is flying a kite on a beach .



GROUND TRUTH CAPTION:

a man riding a big wave on his surfboard
A person surfing under a wave in the ocean.
A surfer catches and rides a large wave.
A person on a surfboard riding a wave.
A person a surf board riding the water waves.

PREDICTED CAPTION:

a man is riding a wave on the water .



Bad Predictions

GROUND TRUTH CAPTION:

A little baby is getting a haircut in a pink chair.
A baby sitting in a chair getting a haircut at a salon.
A barber shop with a young child getting a haircut.
A baby sitting in a chair getting a haircut.
A baby getting a hair cut in the salon sitting in the chair

PREDICTED CAPTION:

a man holding a banana in a banana .



GROUND TRUTH CAPTION:

Two green shoes lined up on a bed.
A pair of sneakers lined up on a bed.
A pair of running shoes sit on the end of a bed.
Bright green sneakers on a bed with a gingham bedspread.
A pair of green sneakers on a single bed with a nightstand next to it

PREDICTED CAPTION:

a bed with a bed and a bed on it .





Thank you!