

# Show, Attend and Tell:

Neural Image Caption Generation With Visual  
Attention (ICML2015)

# “attention”

- Hard attention
- Soft attention
- Two attention based image caption generator model

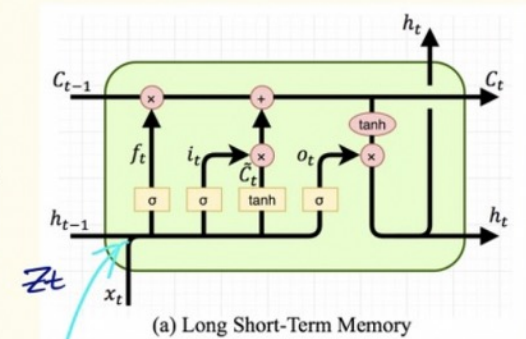
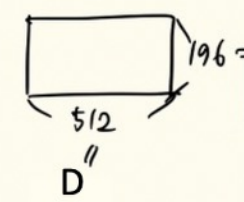
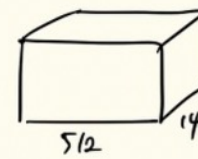
# 1. image captioning model structure

# overview



$14 \times 14 \times 5/2$

$16 \times 5/2$



(a) Long Short-Term Memory

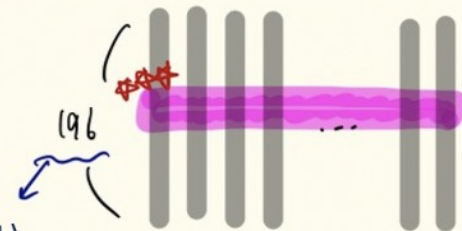
input으로 함께 래입

$z_t$

$$\hat{z}_t = \phi(\{a_i, \alpha_i\})$$

권 생성

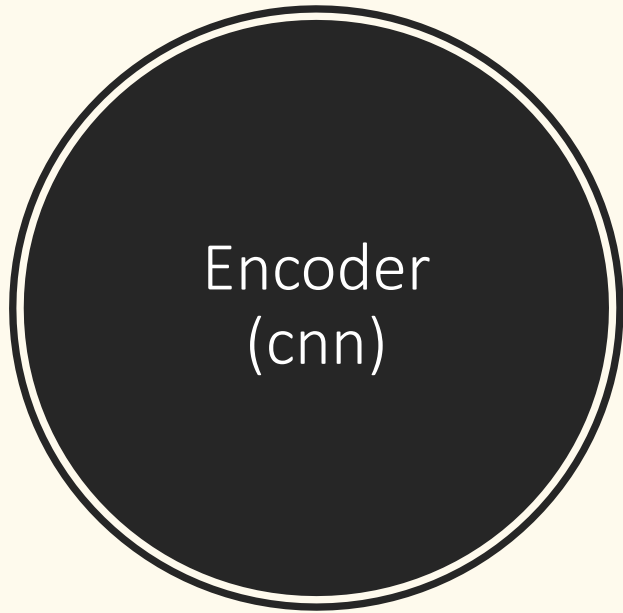
이미지의 지역정보 (L)



특정 위치의 attention!

feature vector의 수 (D):

이미지에서 뽑아낸 특징 5/2개

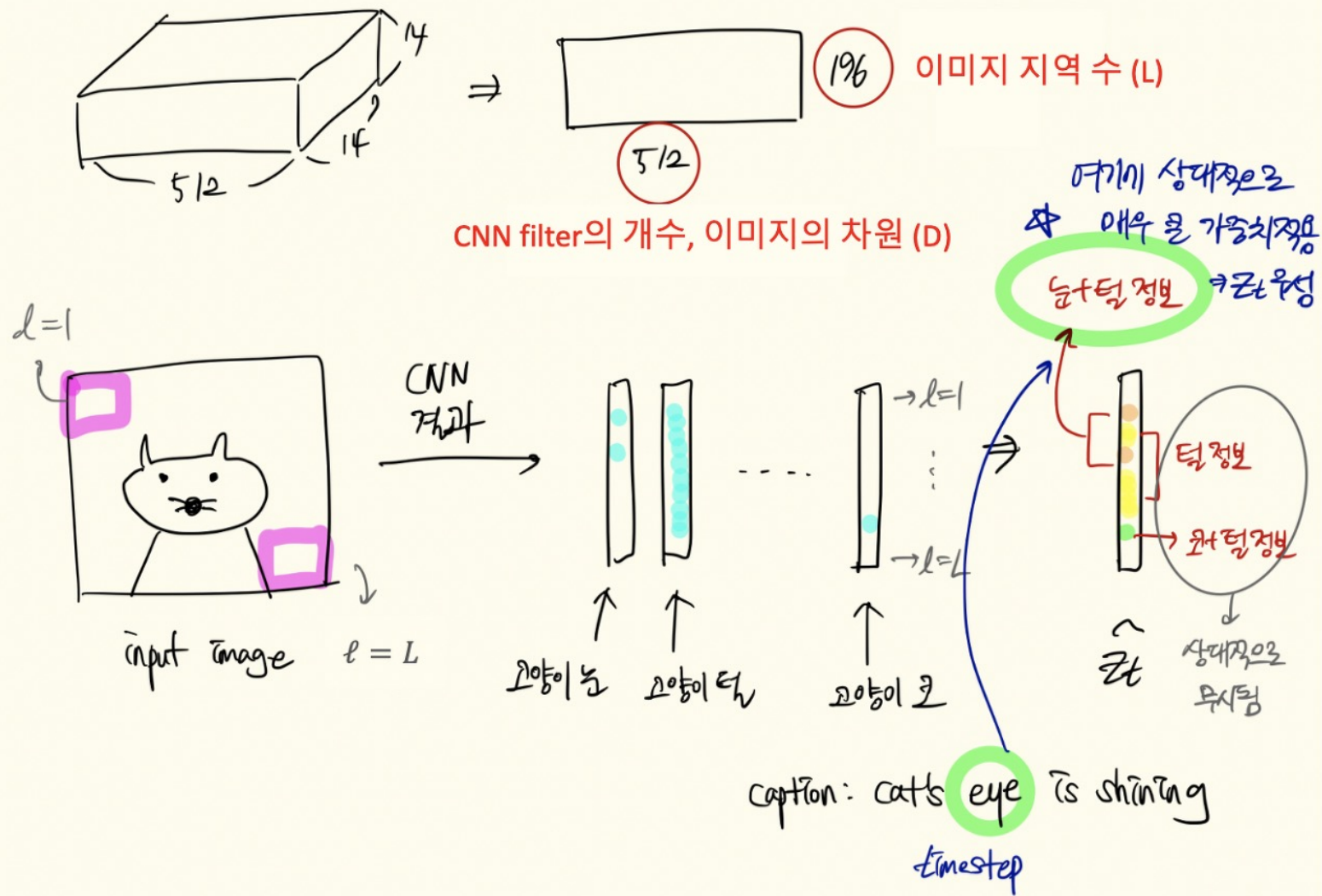


Input: 주어진 이미지

Output: feature vector  $a$  ( $L \times D$ )

- 마지막 layer은 1개의 필터와  $d$ 개의 뉴런
- Fc layer을 사용하지 않음

# Encoder (cnn)



## Decoder (lstm)

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \quad \text{기존 LSTM 수식}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

$$\begin{pmatrix} \mathbf{i}_t \\ \mathbf{f}_t \\ \mathbf{o}_t \\ \mathbf{g}_t \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} T_{D+m+n,n} \begin{pmatrix} \mathbf{E}y_{t-1} \\ \mathbf{h}_{t-1} \\ \hat{\mathbf{z}}_t \end{pmatrix} \quad (1)$$

Attention 적용 후

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{g}_t \quad (2)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t). \quad (3)$$

Input:

1.  $\mathbf{E}y_{t-1}$  : 단어 embedding vector
2.  $\mathbf{h}_{t-1}$  : hidden state vector
3.  $\mathbf{z}_t$  : image feature vector

## 2. Attention Mechanism



$$\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s) = \begin{cases} \mathbf{h}_t^\top \bar{\mathbf{h}}_s & \text{dot} \\ \mathbf{h}_t^\top \mathbf{W}_a \bar{\mathbf{h}}_s & \text{general} \\ \mathbf{v}_a^\top \tanh(\mathbf{W}_a [\mathbf{h}_t; \bar{\mathbf{h}}_s]) & \text{concat} \end{cases}$$

1. Dot (내적): caption의 어느 부분이 이미지의 어느 위치 벡터와 가장 유사한지 파악
2. General: 두 벡터 사이의 연관도를 나타냄
3. Concat: 두 벡터를 concat한 후 학습


### 3. Stochastic Hard Attention

$p(s_{t,i} = 1 \mid s_{j < t}, \mathbf{a}) = \alpha_{t,i} \rightarrow$  timestep  $t$ 까지  $s_{t,i} = 1$ 일 확률  $\alpha_{t,i}$  : multinoulli dist.

$$\hat{\mathbf{z}}_t = \sum_i s_{t,i} \mathbf{a}_i.$$

새로운 R.V (Random Variable)  $\hat{\mathbf{z}}_t$  정의 

$$s_t = [0, 0, 0, 1, \dots, 0, \dots, 0]$$

  
길이 L

# Why stochastic?

- 모든  $s$ 에 대해 계산을 하려면 즉, 모든 위치  $i$ 에 대해 모든 경우의 수를 다 계산하기는 어려움
- 따라서 mini batch와 비슷한 개념으로 sampling하여 gradient를 계산

## 4. Deterministic Soft Attention

$$\mathbb{E}_{p(s_t|a)}[\hat{\mathbf{z}}_t] = \sum_{i=1}^L \alpha_{t,i} \mathbf{a}_i$$

하나만 고르지 않고, 어느 것을 얼마큼 사용할 것인지 비율에 맞게 가져다 사용.  
sampling하지 않아도 되고, end-to-end로 학습이 가능

## 4. 결론



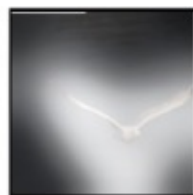
A



bird



flying



over



a



body



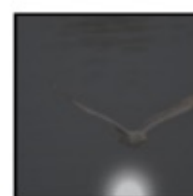
of



water



.





Dataset	Model	BLEU				METEOR
		BLEU-1	BLEU-2	BLEU-3	BLEU-4	
Flickr8k	Google NIC(Vinyals et al., 2014) <sup>†Σ</sup>	63	41	27	—	—
	Log Bilinear (Kiros et al., 2014a) <sup>◦</sup>	65.6	42.4	27.7	17.7	17.31
	Soft-Attention	<b>67</b>	44.8	29.9	19.5	18.93
	Hard-Attention	<b>67</b>	<b>45.7</b>	<b>31.4</b>	<b>21.3</b>	<b>20.30</b>
Flickr30k	Google NIC <sup>†◦Σ</sup>	66.3	42.3	27.7	18.3	—
	Log Bilinear	60.0	38	25.4	17.1	16.88
	Soft-Attention	66.7	43.4	28.8	19.1	<b>18.49</b>
	Hard-Attention	<b>66.9</b>	<b>43.9</b>	<b>29.6</b>	<b>19.9</b>	18.46
COCO	CMU/MS Research (Chen & Zitnick, 2014) <sup>a</sup>	—	—	—	—	20.41
	MS Research (Fang et al., 2014) <sup>†a</sup>	—	—	—	—	20.71
	BRNN (Karpathy & Li, 2014) <sup>◦</sup>	64.2	45.1	30.4	20.3	—
	Google NIC <sup>†◦Σ</sup>	66.6	46.1	32.9	24.6	—
	Log Bilinear <sup>◦</sup>	70.8	48.9	34.4	24.3	20.03
	Soft-Attention	70.7	49.2	34.4	24.3	<b>23.90</b>
	Hard-Attention	<b>71.8</b>	<b>50.4</b>	<b>35.7</b>	<b>25.0</b>	23.04

- **BLEU score**: N-gram 단위로 정답 caption과의 유사도를 측정
- Flickr dataset 및 COCO dataset 모두에서 **attention**을 적용한 모델의 성능(BLEU score)이 더 높음