# Neural Image Caption Generation with Visual Attention

Kelvin Xu, Jimmy Lei Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard S. Zemel, Yoshua Bengio

Literature Review
Deep Learning - DPL302m

Today's Date, 2024

# Paper Overview

- Introduction
- Related Work
- Image Caption Generation with Attention Mechanism
    - Encoder: Convolutional Features
    - Decoder: LSTM Network
- Learning Stochastic "Hard" vs Deterministic "Soft" Attention
    - Stochastic "Hard" Attention
    - Deterministic "Soft" Attention
- Training & Evaluation
- Conclusion

# Content Overview

- Introduction: Why Attention?

- Challenges

- Model
    - Architecture

    - Attention Mechanism
        - Soft Attention
        - Hard Attention

# Introduction



Visual System
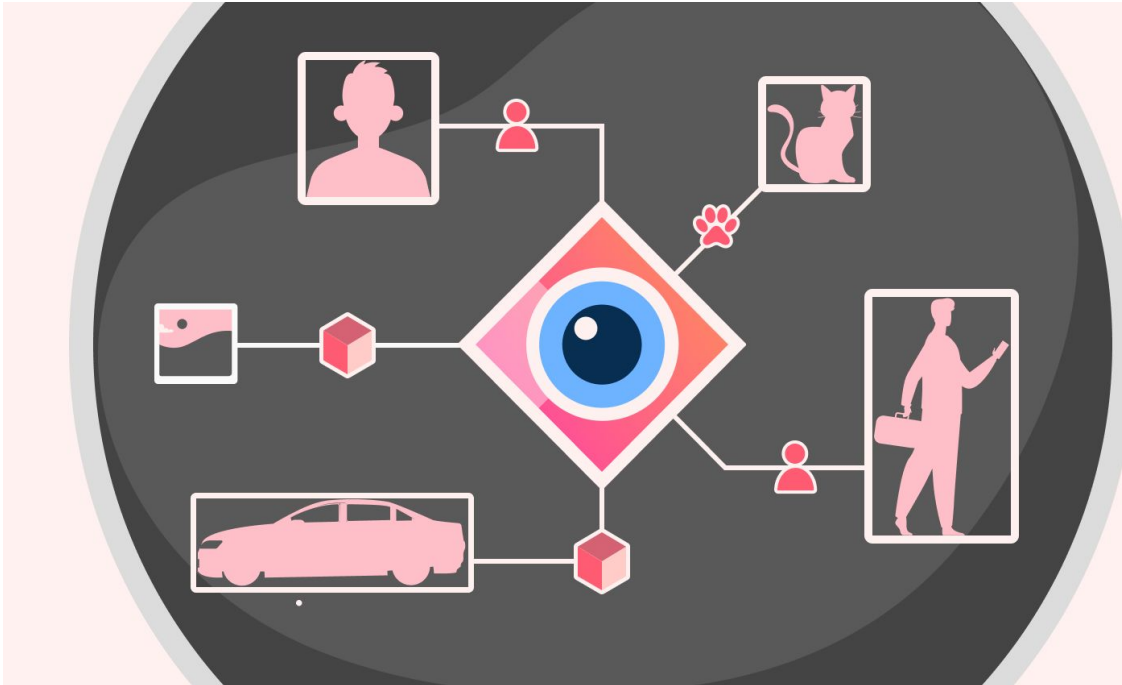


Language System

# Introduction

# Introduction

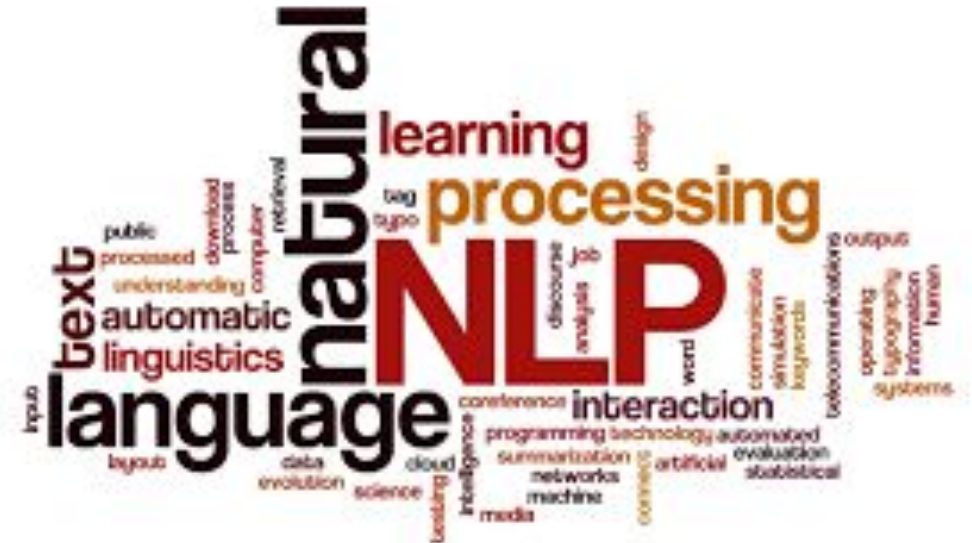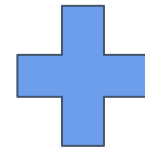Scene Understanding $\rightarrow$ Better Evaluation for Computer Vision System

# Introduction



Computer Vision

+

Natural Language Processing
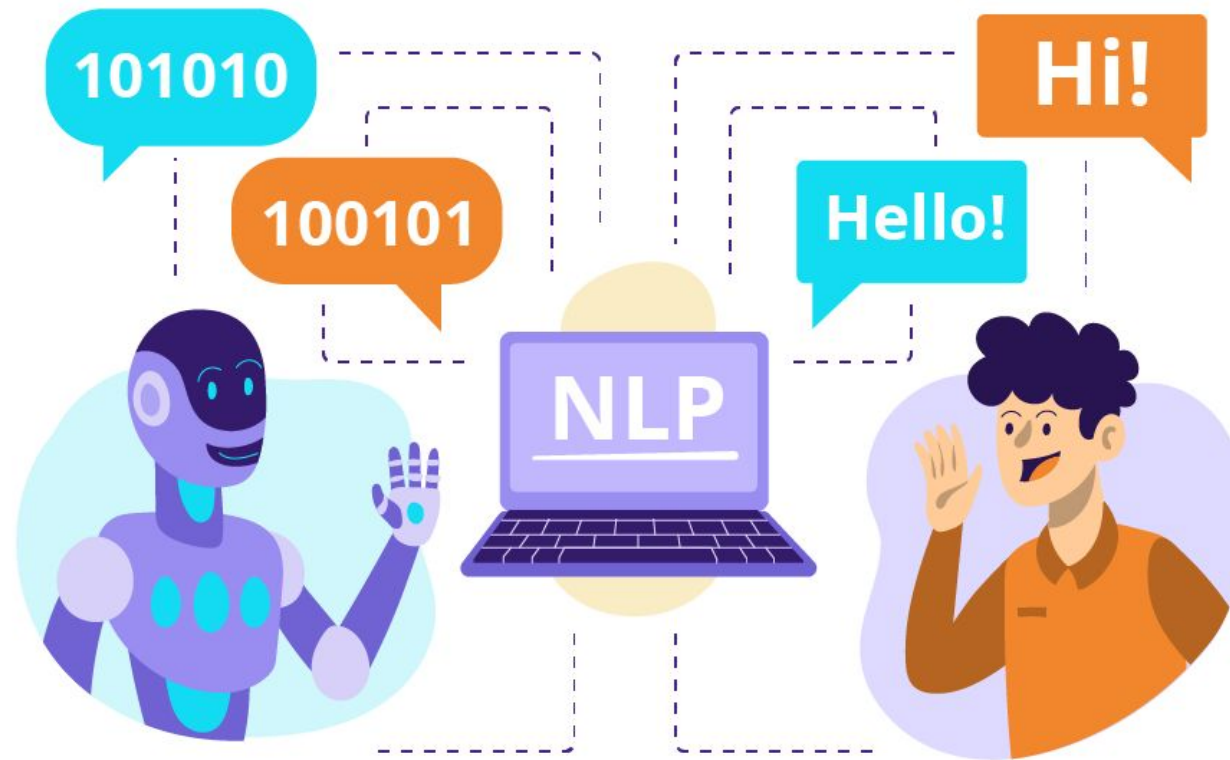
# Introduction

Object Recognition

The most IMPORTANT object

Relationship between objects

# Introduction

Practical Meaning: AI can finally understand daily scenes and INTERACT with human!

# Introduction
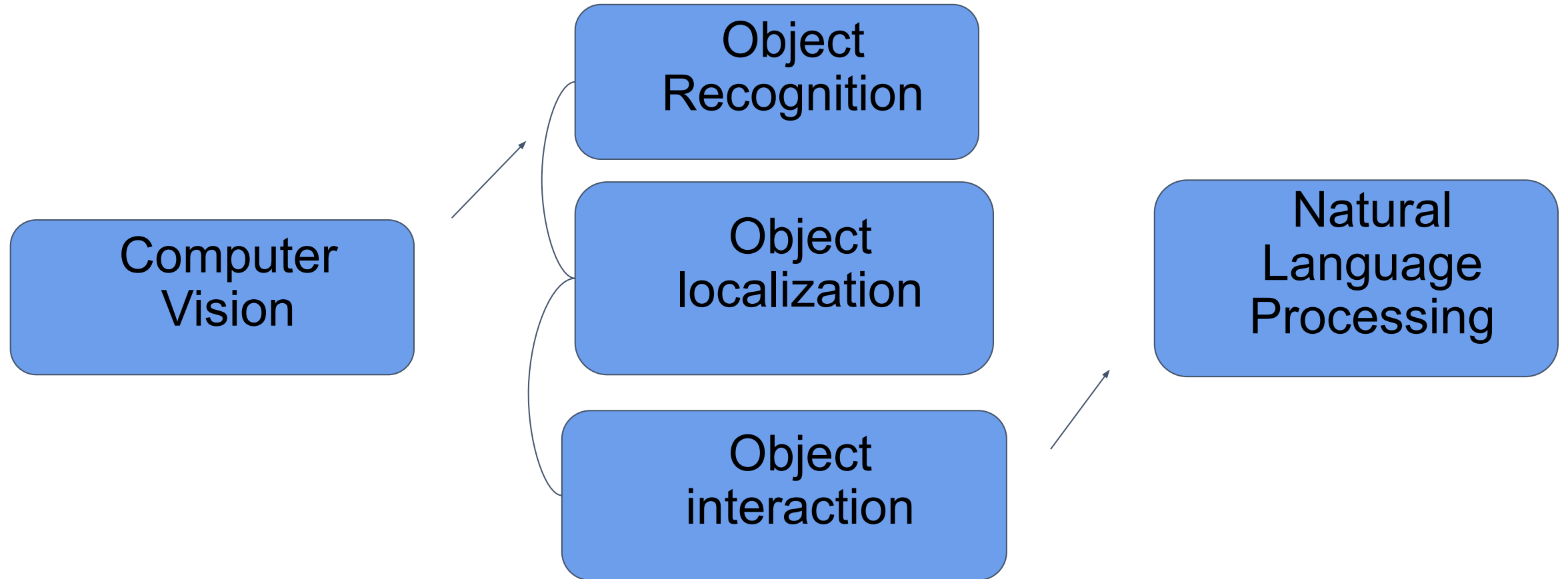
But Daily scenes ⇒ Visual Noises!

# Challenges

- How to extract visual information from the image

- How to transform the visual information into proper natural language

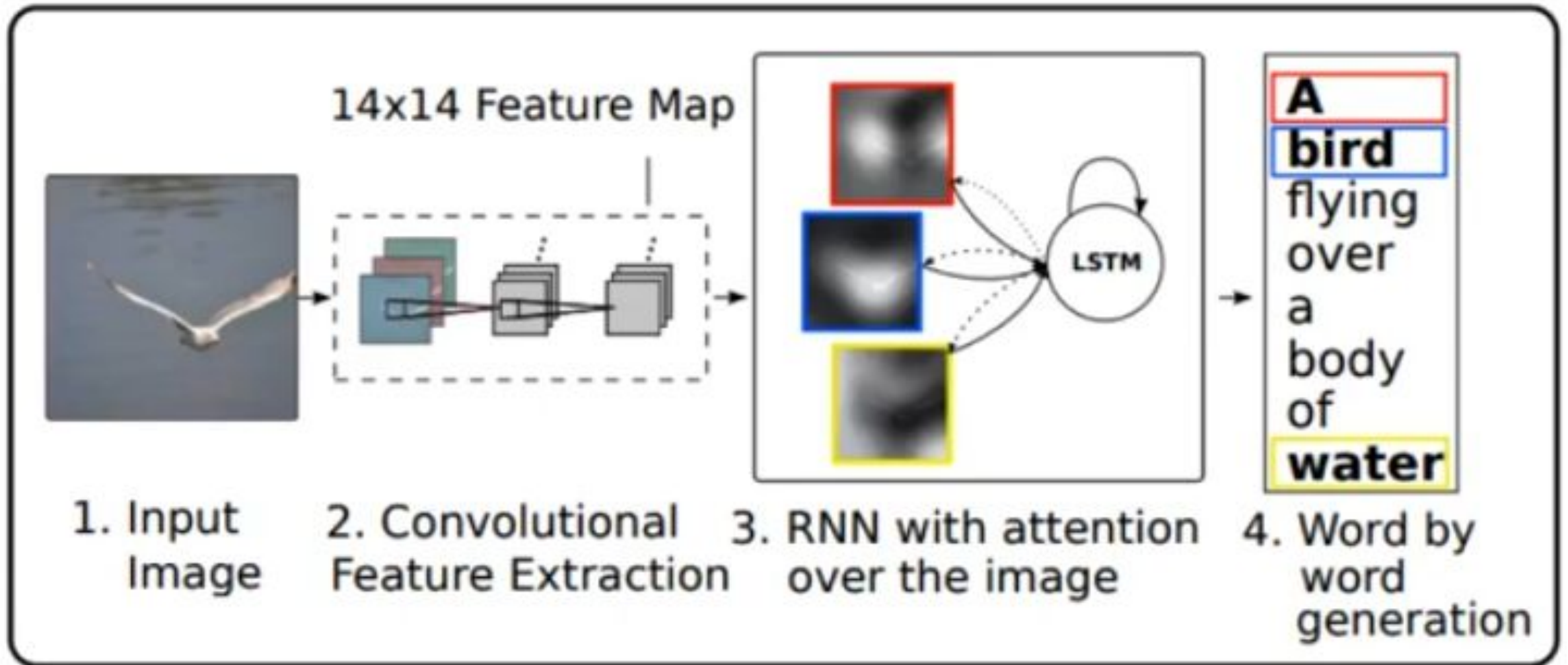- **How to represent different parts of images and let the model focus on the IMPORTANT part**

# Model

Divide the image caption generator into 2 sub-problems:

- Image understanding

- Descriptive caption generating

# Introduction

# Architecture



1. Input Image
2. Convolutional Feature Extraction
3. RNN with attention over the image
4. Word by word generation

# Convolutional Neural Network (CNN)

CNN consists of three types of layers:

- Convolutional layer

- Pooling layer

- Fully Connected layer

In this paper, they use a pre-trained CNN model VGG-Net 16 to do the image understanding.

# NLP Part - Dependency between output

In NLP part, each output is a word in the caption. Usually the words in the one caption are highly related to each other



*A girl holding her teddy bear*

# Long-short Term Memory Network (LSTM)

**Problem:** The output produced by most neural networks are independent
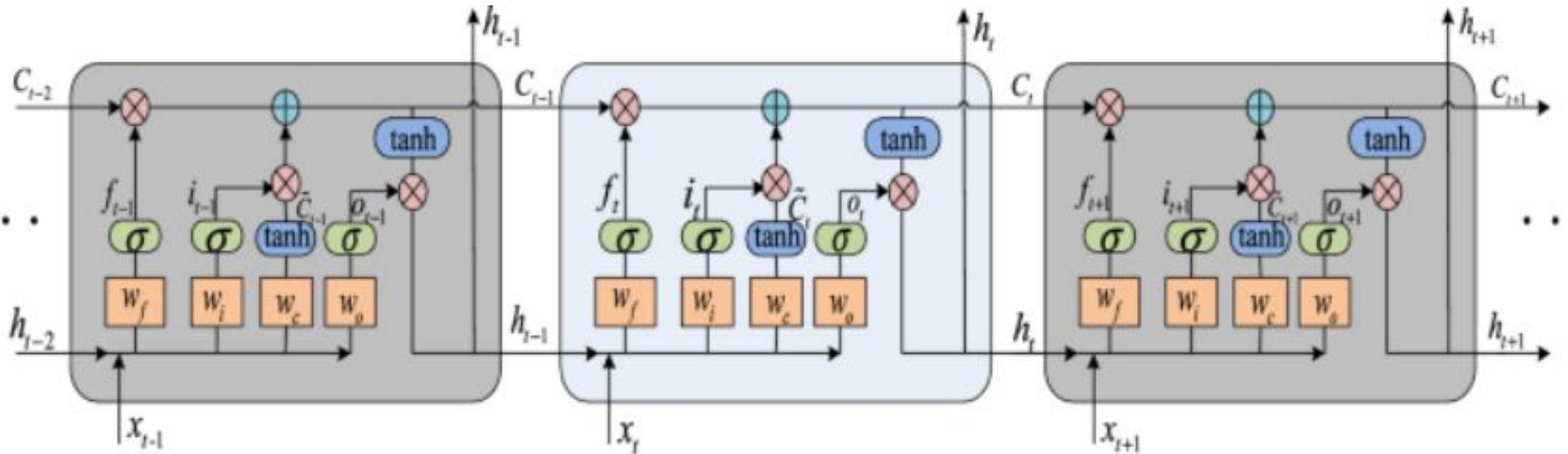
Solution:

Long-short Term memory network (LSTM)



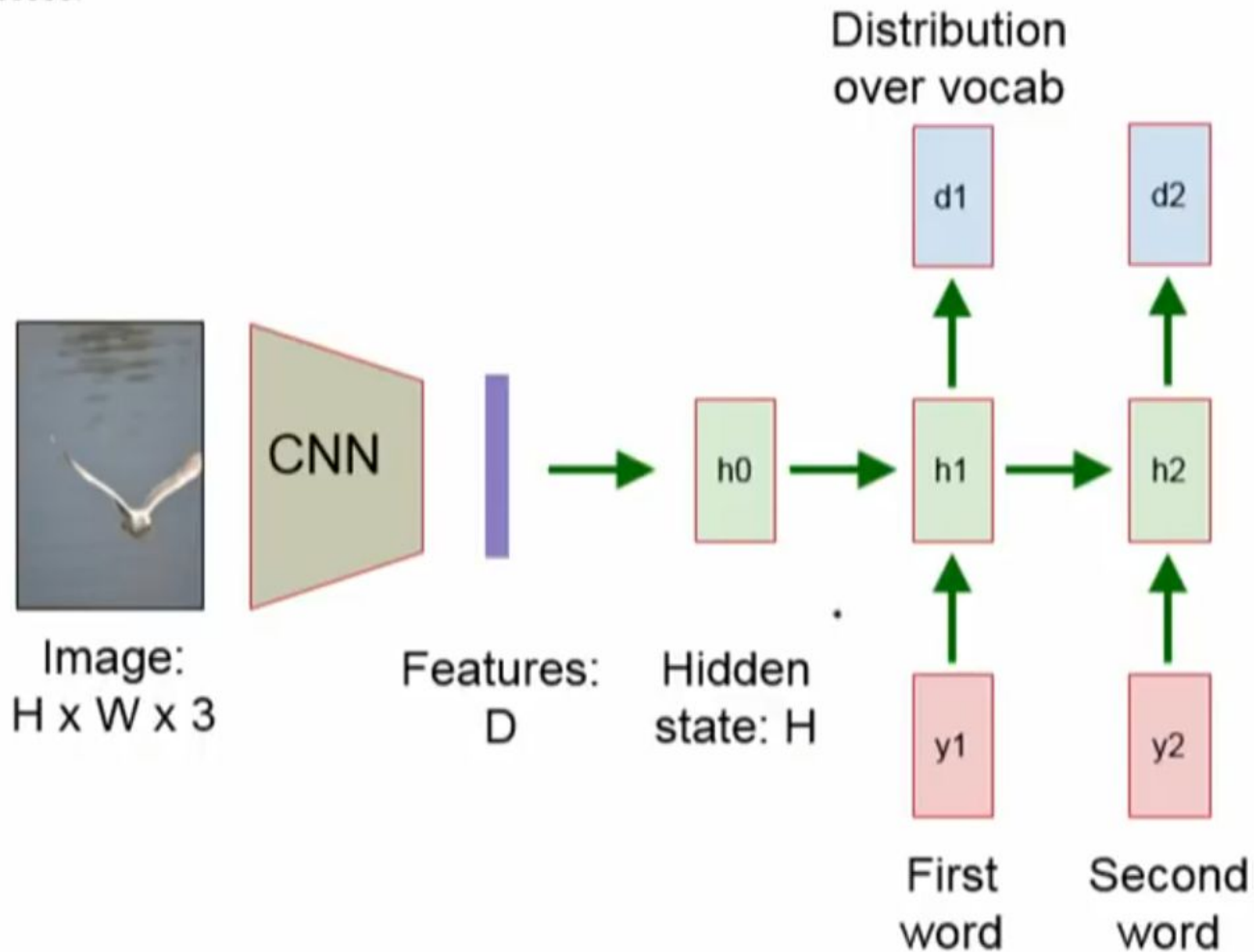*A <u>girl</u> holding <u>her</u> teddy bear*

# Long-short Term Memory Network (LSTM)



LSTM has a powerful memory mechanism.

LSTM takes a new input and also inherits the hidden state and memory cell of the previous step when producing a new output.
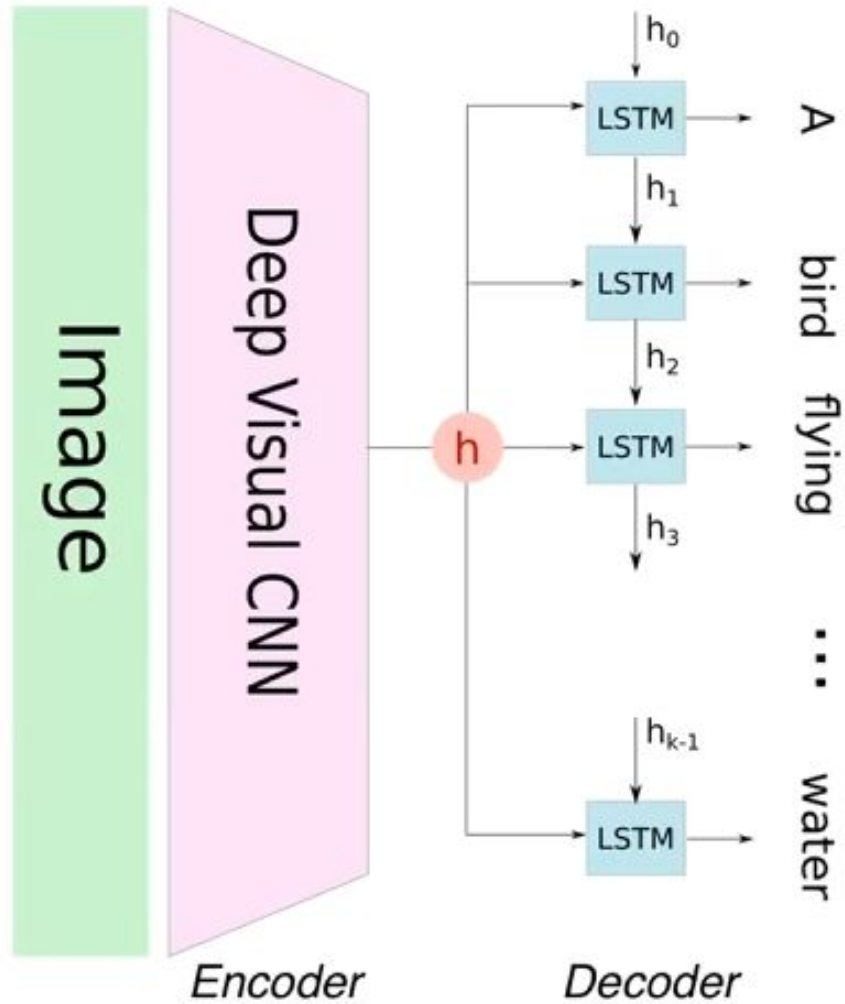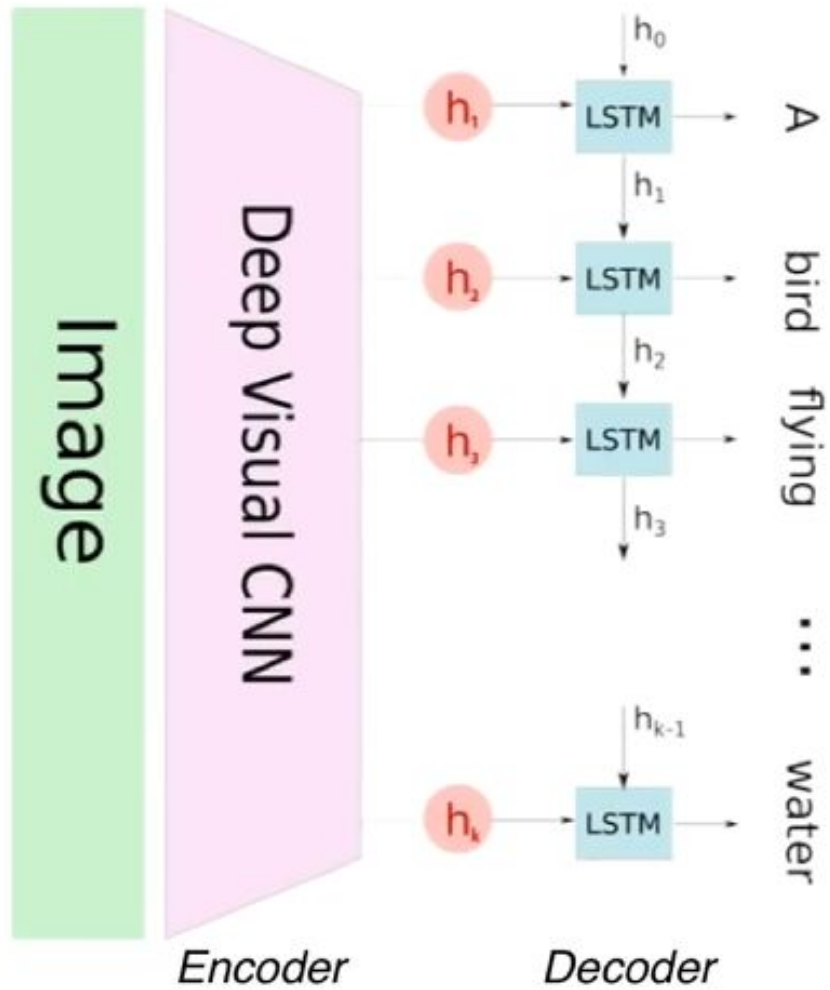
# Attention Mechanism



Distribution over vocab

Image: H x W x 3

CNN

Features: D

Hidden state: H

First word

Second word

▶ RNN looks at the whole image only once

▶ What if RNN looks at different parts of the image at each time step?

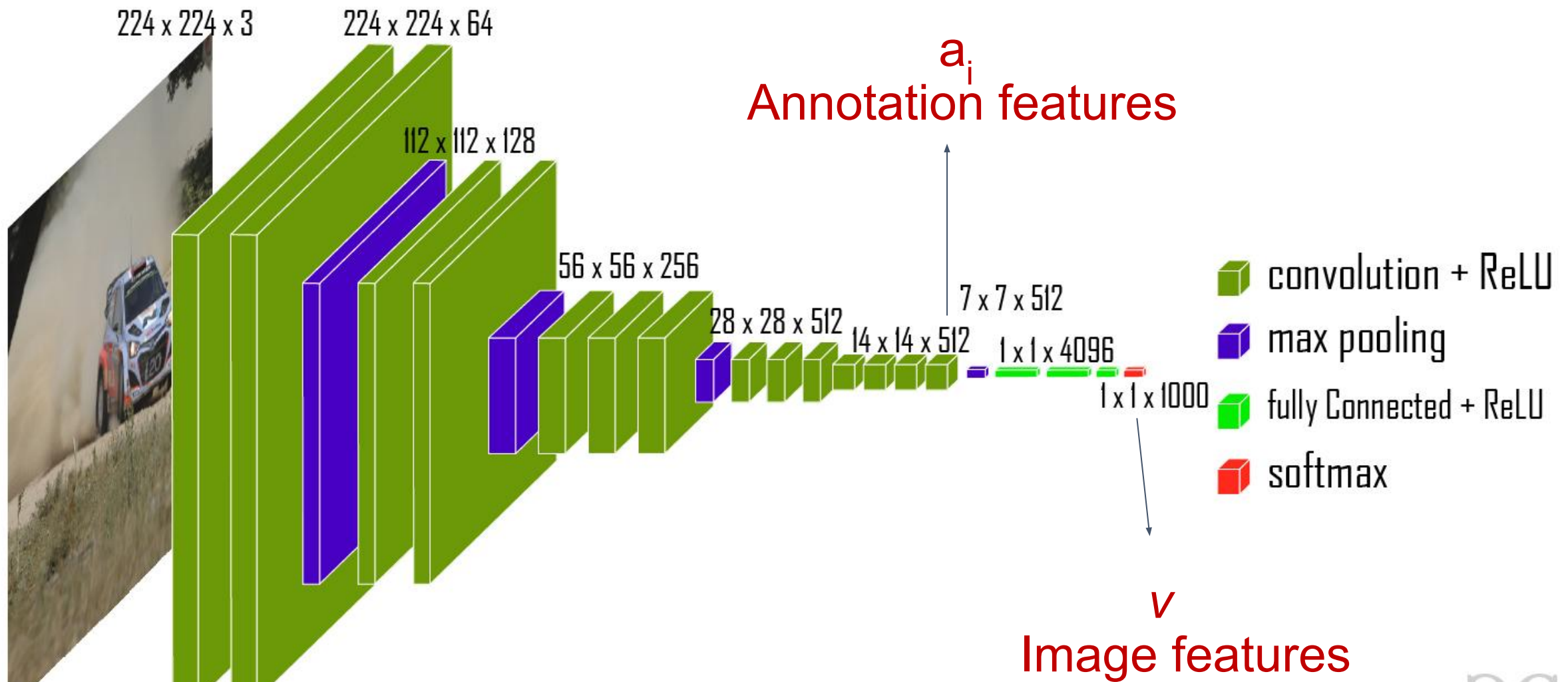# Model without Attention

# Model with Attention

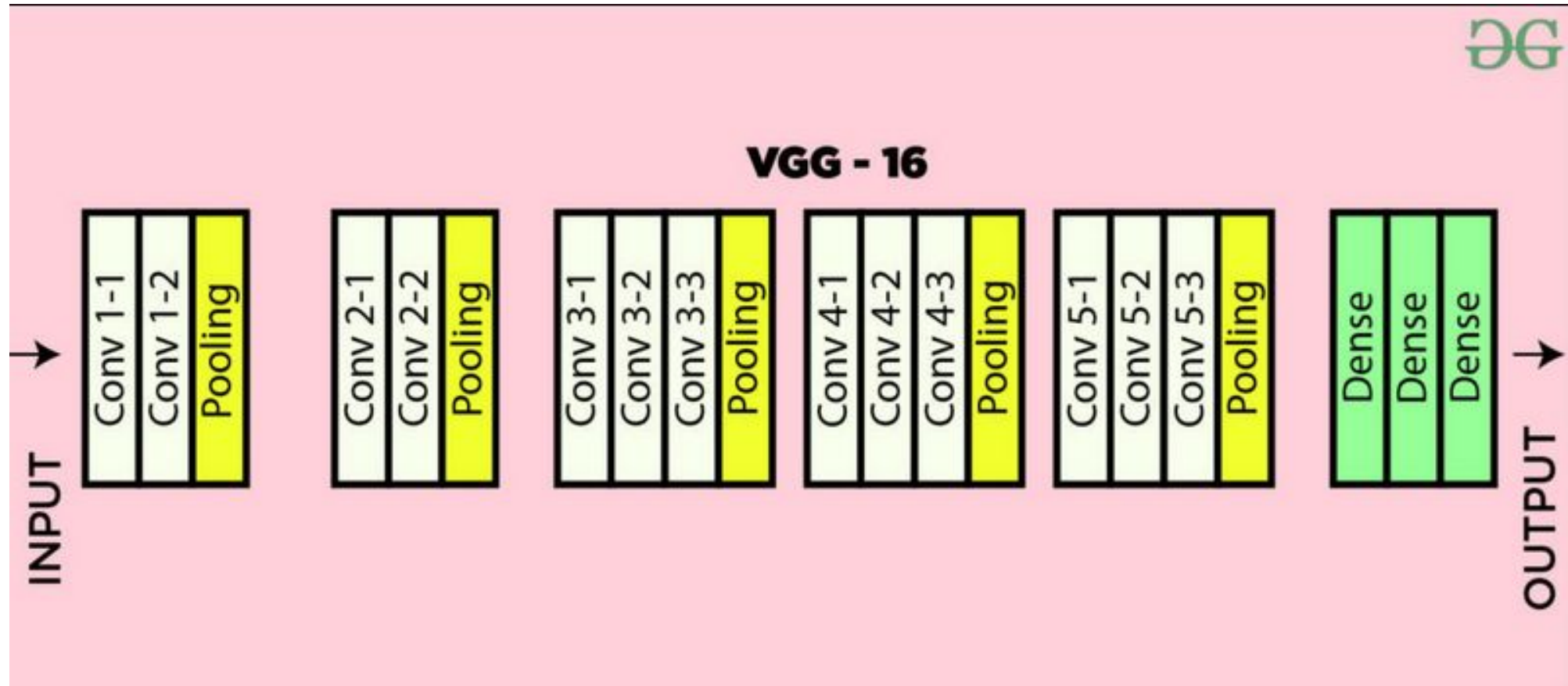# Model with Attention

### 3.1.1. ENCODER: CONVOLUTIONAL FEATURES

We use a convolutional neural network in order to extract a set of feature vectors which we refer to as annotation vectors. The extractor produces $L$ vectors, each of which is a D-dimensional representation corresponding to a part of the image.

$$a = \{\mathbf{a}_1, \ldots, \mathbf{a}_L\}, \, \mathbf{a}_i \in \mathbb{R}^D$$

224 x 224 x 3
224 x 224 x 64
112 x 112 x 128
56 x 56 x 256
28 x 28 x 512
14 x 14 x 512
7 x 7 x 512
1 x 1 x 4096
1 x 1 x 1000

$a_i$
Annotation features

$v$
Image features

convolution + ReLU
max pooling
fully Connected + ReLU
softmax

# Annotation features

# Annotation features



Conv 5-3 Layer (subset)

# Attention Part



CNN (2) → Annotation features $\{a_1, a_2, ..., a_n\}$

LSTM → Hidden state $h_{t-1}$

Annotation features and Hidden state → Attention Mechanism → Attention Vector $z_t$

# Model with attention

# Attention Mechanism



Image: H x W x 3

CNN

Features: L x D

# Attention Mechanism



Distribution over
L locations

a1

h0

Image:
H x W x 3

CNN

Features:
L x D

# Attention Mechanism



Distribution over L locations

CNN

a1

h0

Image: H x W x 3

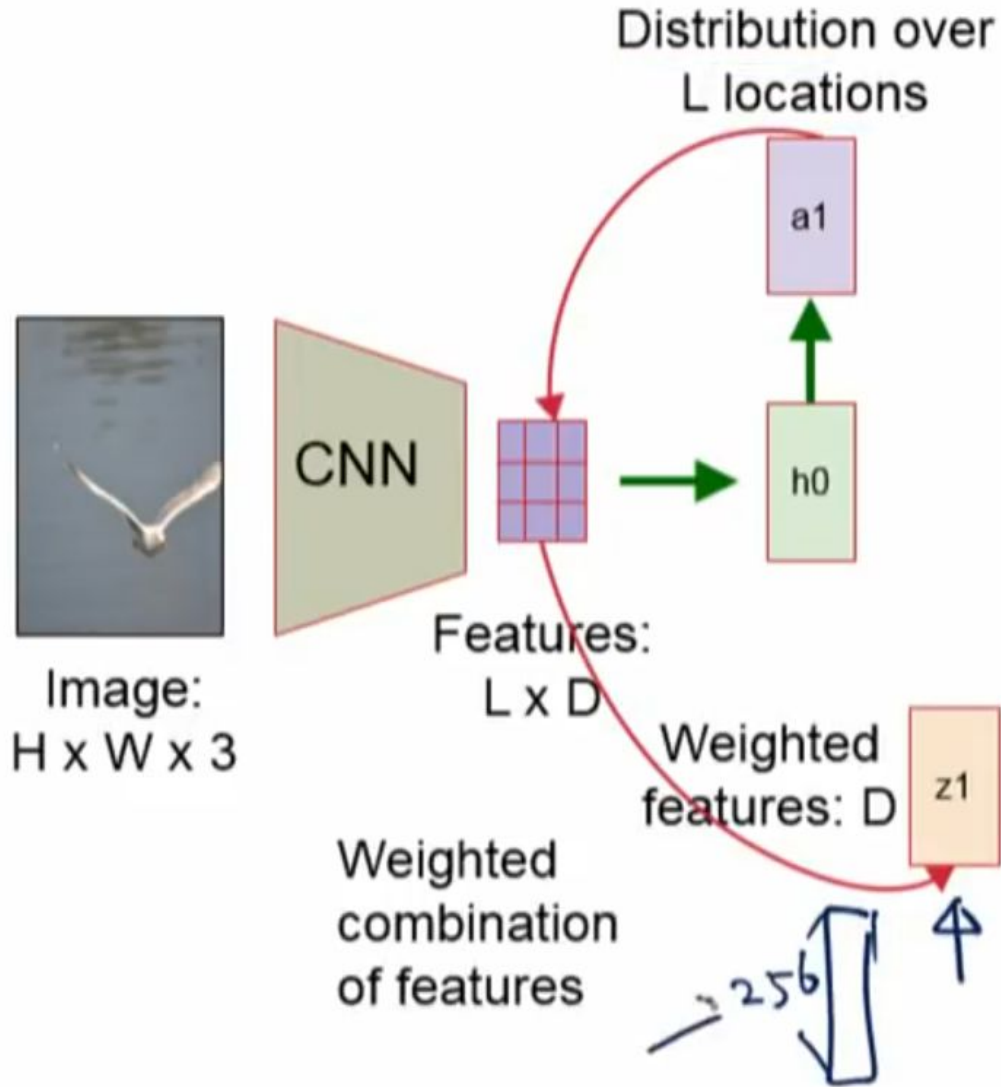Features: L x D

Weighted features: D    z1
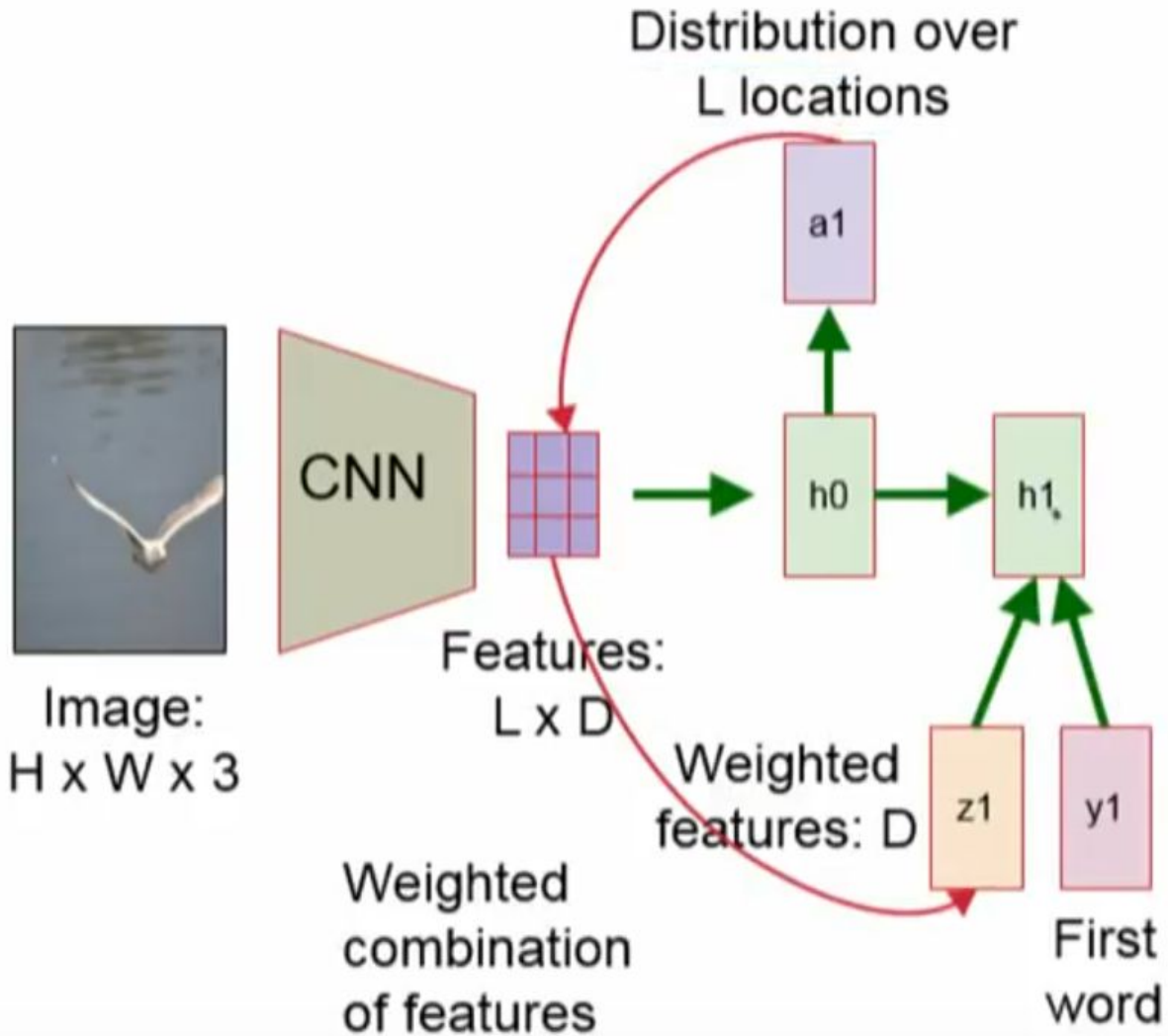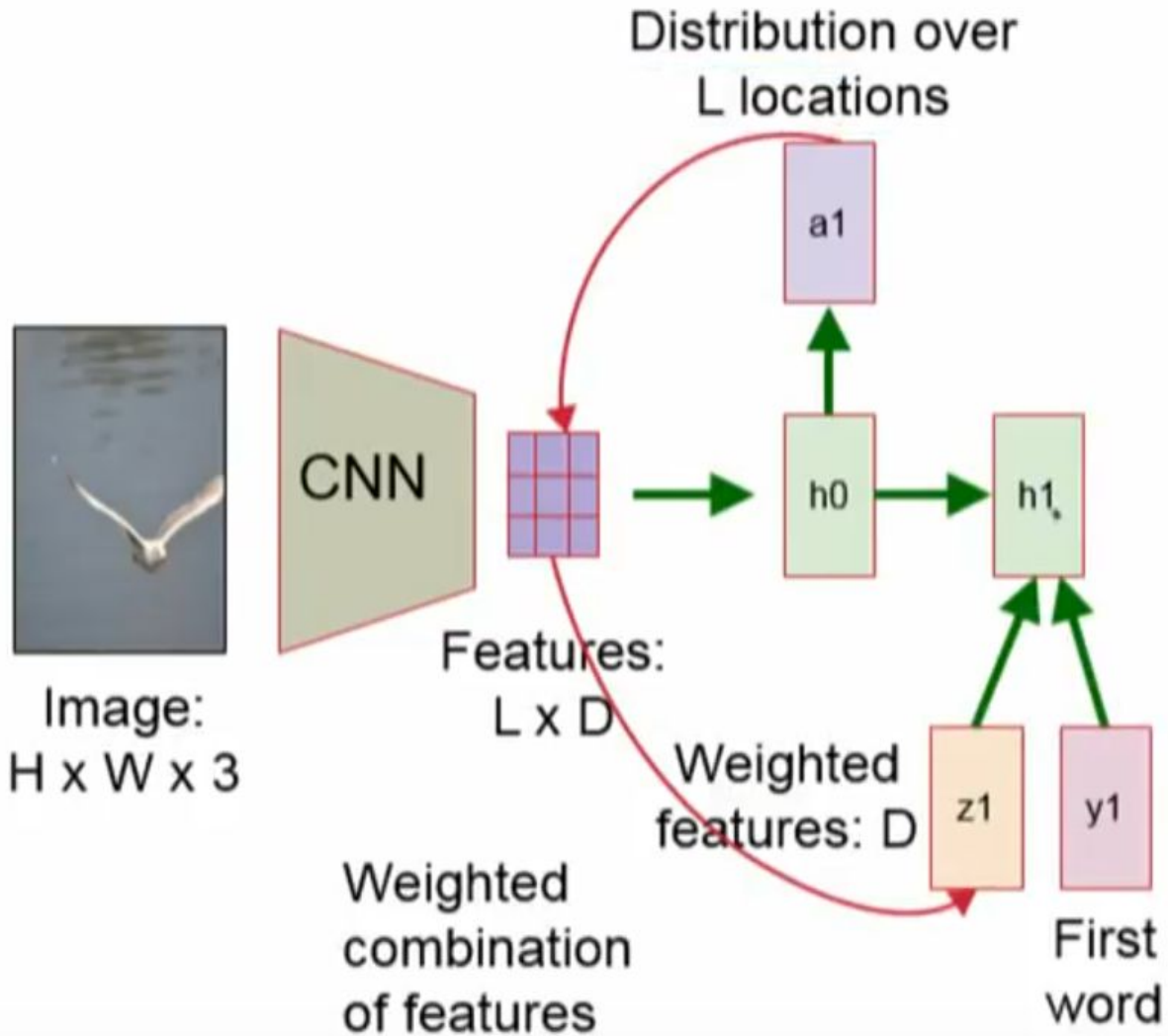
Weighted combination of features

256

# Attention Mechanism

# Attention Mechanism

# Attention Mechanism

# Attention Mechanism

# Attention Unit



$$m_i = \tanh\left(y_i W_{y_i} + C W_c\right)$$

# Attention Unit

# Attention Unit



$$\hat{\mathbf{z}}_t = \phi\left(\{\mathbf{a}_i\}, \{\alpha_i\}\right) = \sum_i \alpha_i \, \mathbf{a}_i$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^{L} \exp(e_{tk})}.$$

$$e_{ti} = f_{\text{att}}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$= \tanh(W_c C + W_x X_i) = \tanh(W_c h_{t-1} + W_x x_i)$$

- **Soft Attention**: different parts, different subregions

- **Hard Attention**: only ONE subregion

# Type of Attention

- **Soft Attention**: different parts, different subregions

$$z = \sum_n s_n y_n$$

- Soft Attention is <u>Deterministic</u>





$a = $ "Move 5ft. Forward"

ALWAYS

$s=\{(0,0), \text{Forward}\}$         $s'=\{(5,0), \text{Forward}\}$

# Type of Attention

- **Soft Attention**: different parts, different subregions

$$z = \sum_n s_n y_n$$

- Soft Attention is <u>Deterministic</u>



Soft Attention →

# Implementing Soft Attention



CNN

Image:
H x W x 3

Grid of features
(Each
D-dimensional)

$$e_{ti} = f_{\text{att}}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

Context vector
z
(D-dimensional)

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^{L} \exp(e_{tk})}.$$

From
RNN:

Distribution over
grid locations
$p_a + p_b + p_c + p_c = 1$

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

(b) A stop sign is on a road with a mountain in the background.

(b) A woman holding a clock in her hand.

# Type of Attention

## 2. Hard Attention: only ONE subregion

Hard Attention is *Stochastic*

# Implementing Hard Attention

$$p(s_{t,i} = 1 \mid s_{j<t}, \mathbf{a}) = \alpha_{t,i}$$

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

We represent the location variable $s_t$ as where the model decides to focus attention when generating the $t^{th}$ word. $s_{t,i}$ is an indicator one-hot variable which is set to 1 if the $i$-th location (out of $L$) is the one used to extract visual features. By treating the attention locations as intermediate latent variables, we can assign a multinoulli distribution parametrized by $\{\alpha_i\}$, and view $\hat{z}_t$ as a random variable:
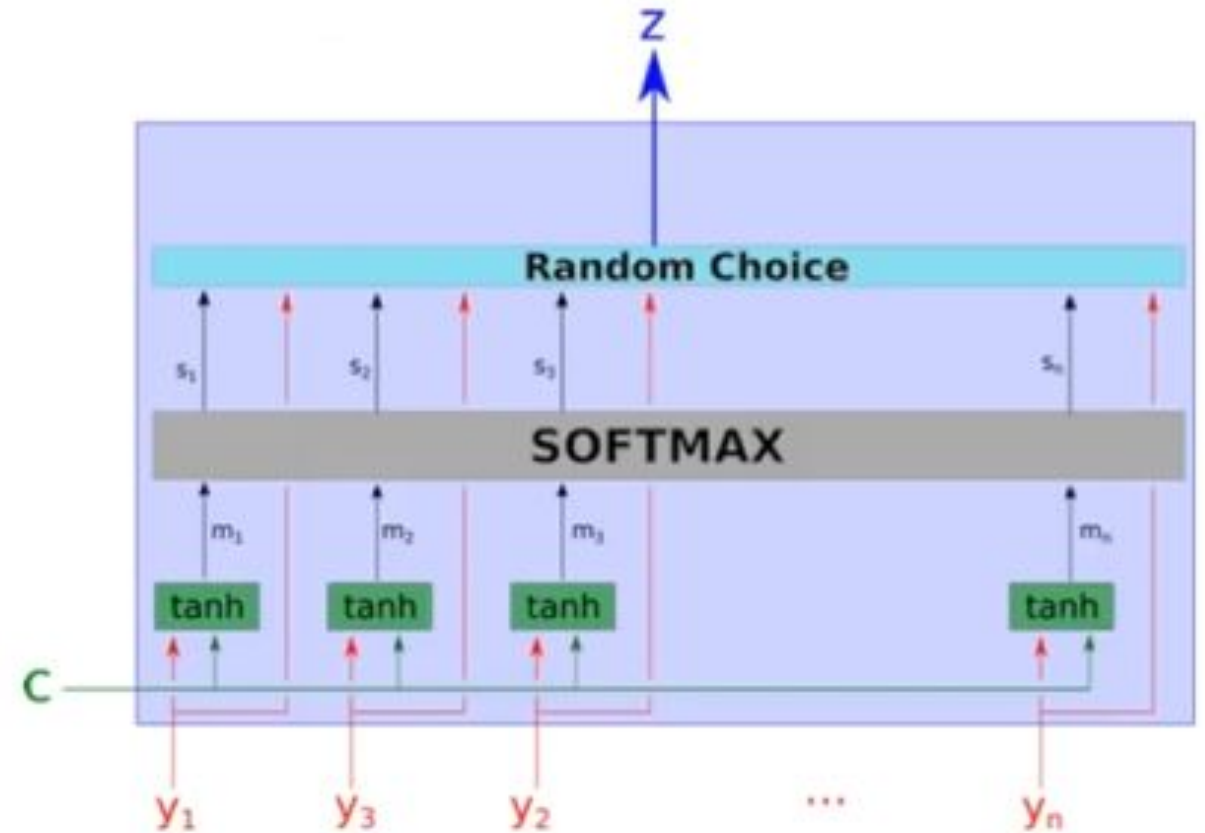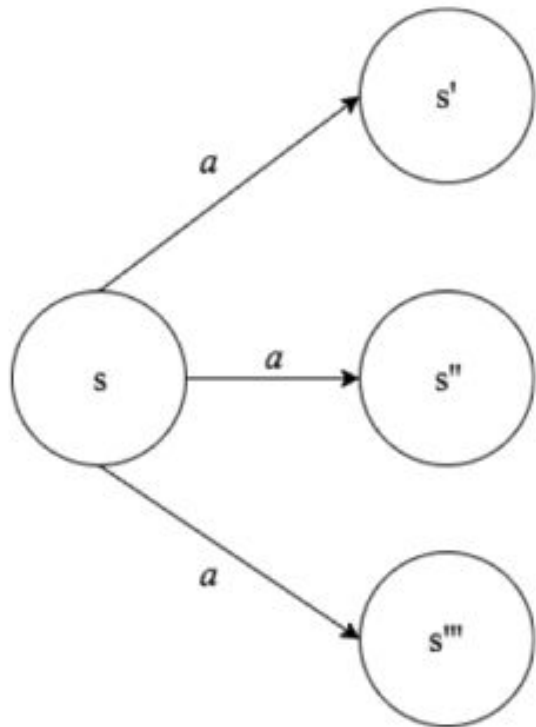
$$L_s = \sum_s p(s \mid \mathbf{a}) \log p(\mathbf{y} \mid s, \mathbf{a})$$

$$\leq \log \sum_s p(s \mid \mathbf{a}) p(\mathbf{y} \mid s, \mathbf{a})$$

$$= \log p(\mathbf{y} \mid \mathbf{a})$$

We define a new objective function $L_s$ that is a variational lower bound on the marginal log-likelihood $\log p(\mathbf{y} \mid \mathbf{a})$ of observing the sequence of words y given image features a. The learning algorithm for the parameters $W$ of the models can be derived by directly optimizing $L_s$:

# Implementing Hard Attention

$$\frac{\partial L_s}{\partial W} = \sum_s p(s \mid \mathbf{a}) \left[ \frac{\partial \log p(\mathbf{y} \mid s, \mathbf{a})}{\partial W} + \right.$$

$$\left. \log p(\mathbf{y} \mid s, \mathbf{a}) \frac{\partial \log p(s \mid \mathbf{a})}{\partial W} \right]. \quad (11)$$

$$= \sum_s p(s \mid \mathbf{a}) \frac{\partial \log p(\mathbf{y} \mid s, \mathbf{a})}{\partial W} + p(s \mid \mathbf{a}) \log p(\mathbf{y} \mid s, \mathbf{a}) \frac{\partial \log p(s \mid \mathbf{a})}{\partial W}$$

# Implementing Hard Attention

- This means that Monte Carlo Sampling can be performed!

$$\tilde{s}_t \sim \text{Multinoulli}_L(\{\alpha_i\})$$

$$\frac{\partial L_s}{\partial W} \approx \frac{1}{N} \sum_{n=1}^{N} \left[ \frac{\partial \log p(\mathbf{y} \mid \tilde{s}^n, \mathbf{a})}{\partial W} + \log p(\mathbf{y} \mid \tilde{s}^n, \mathbf{a}) \frac{\partial \log p(\tilde{s}^n \mid \mathbf{a})}{\partial W} \right]$$

# Implementing Hard Attention

- The issue is that the variance in this estimate is too high

$$b_k = 0.9 \times b_{k-1} + 0.1 \times \log p(\mathbf{y} \mid \tilde{s}_k, \mathbf{a})$$

$$\frac{\partial L_s}{\partial W} \approx \frac{1}{N} \sum_{n=1}^{N} \left[ \frac{\partial \log p(\mathbf{y} \mid \tilde{s}^n, \mathbf{a})}{\partial W} + \lambda_r (\log p(\mathbf{y} \mid \tilde{s}^n, \mathbf{a}) - b) \frac{\partial \log p(\tilde{s}^n \mid \mathbf{a})}{\partial W} + \lambda_e \frac{\partial H[\tilde{s}^n]}{\partial W} \right]$$
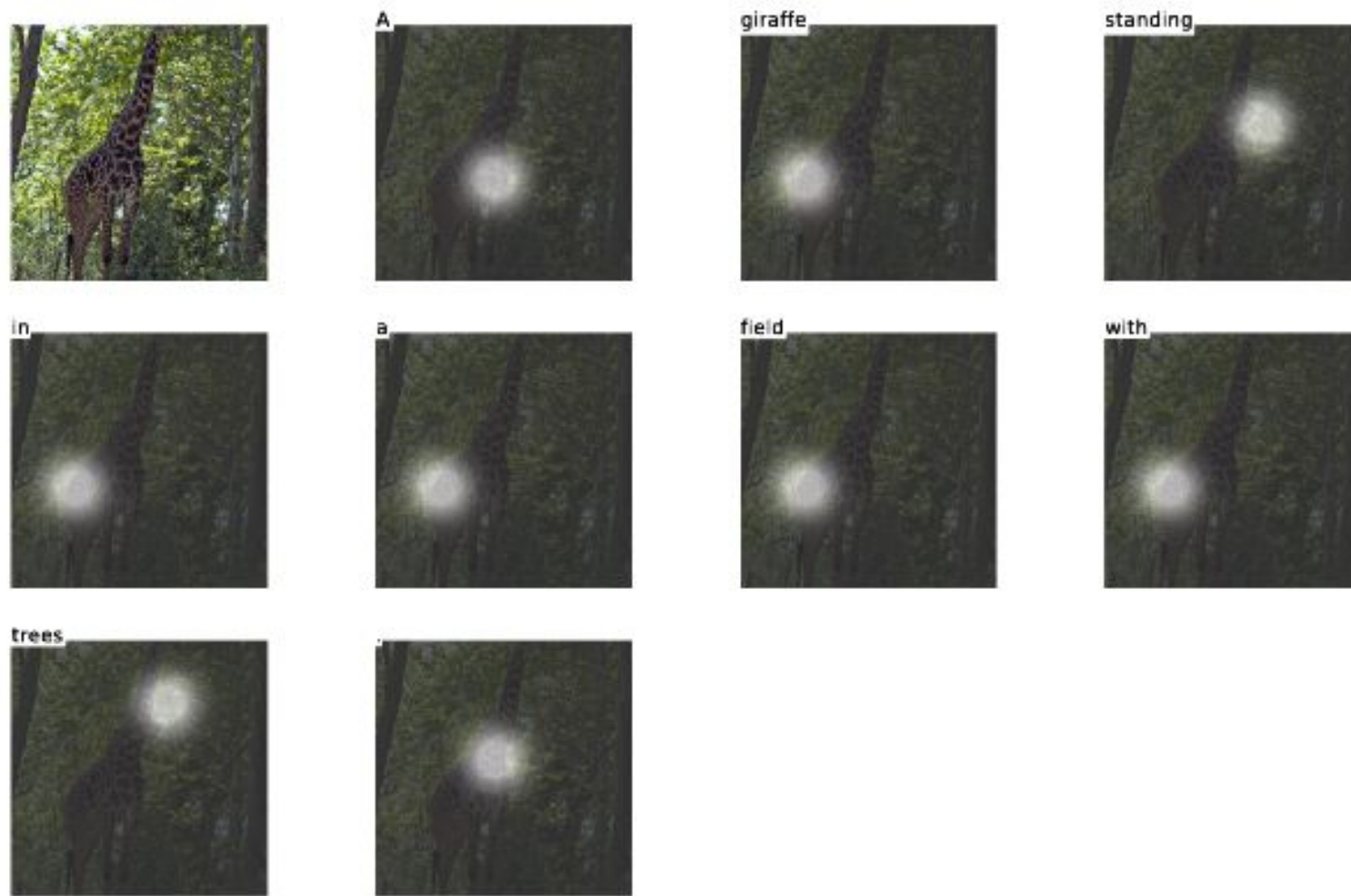
(a) A giraffe standing in a field with trees.

# Negative Example



(a) A dog is laying on a bed with a book.

# Doubly Stochastic Attention

- To encourage the model to look at various parts of the image

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_i^L (1 - \sum_t^C \alpha_{ti})^2$$

# Result

Table 1. BLEU-1,2,3,4/METEOR metrics compared to other methods, † indicates a different split, (—) indicates an unknown metric, ○ indicates the authors kindly provided missing metrics by personal communication, Σ indicates an ensemble, $a$ indicates using AlexNet

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