

# Introduction to Attention Models

**Dr. Dileep A. D.**

Associate Professor,  
Multimedia Analytics Networks And Systems (MANAS) Lab,  
School of Computing and Electrical Engineering (SCEE),  
Indian Institute of Technology Mandi, Kamand, H.P.  
Email: [addileep@iitmandi.ac.in](mailto:addileep@iitmandi.ac.in)



## Agenda

- Neural network models
  - Fully connected neural networks
  - Convolutional neural networks
  - Recurrent neural networks
- Introduction to encoder-decoder models
- Attention Mechanism
  - Attention models in vision

## Artificial Neural Networks

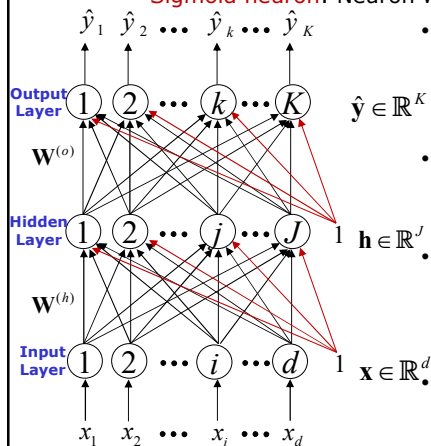
- **Learning method:**
  - Error correction learning (Backpropagation algorithm [1])
- **Structure of network:**
  - **Feedforward neural networks**
    - Fully Connected Neural Network (FCNN),
    - Convolutional Neural Networks (CNN),
    - Auto Encoders
  - **Feedback neural networks**
    - Recurrent Neural Networks (RNN)
    - Long Short Term Memory (LSTM)
  - **Feedforward and feedback neural networks**
    - Bidirectional LSTM
    - Self Organizing Maps (SOM)

[1] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Learning internal representations by error propagation. In D. E. Rumelhart and J. L. McClelland, editors, *Parallel Distributed Processing*, volume 1, pages 318-362. MIT Press, 1986.

3

## Fully Connected Neural Network (FCNN)

- Architecture of an FCNN:
- **Input layer:** Linear neurons
  - **Linear neuron:** When an input is given to a neuron and the same input comes out as output
- **Hidden layers (1 or 2 or more):** Sigmoidal neurons or ReLU
  - **Sigmoid neuron:** Neuron with sigmoid activation function



- **Output layer:** **Sigmoidal/Softmax neurons** (for pattern classification task) or **Linear neuron** (for regression)

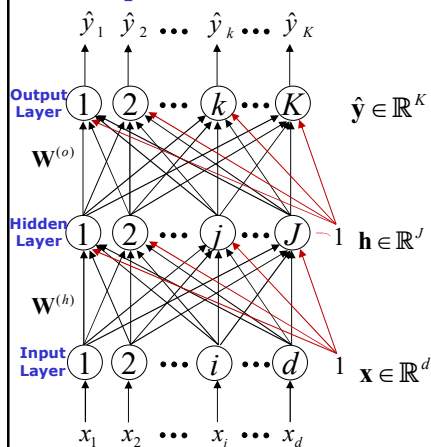
- Number of neurons in input layer ( $d$ ): **Dimension of the data vector** (number of input variables)

- Number of neurons in the output layer ( $K$ ): **Number of classes in classification or number of output variables**

- Number of layers and neurons in each of the hidden layers are decided experimentally

4

## Fully Connected Neural Network (FCNN)



$$\mathbf{h} = g(\mathbf{W}^{(h)\top} \mathbf{x} + \mathbf{w}_0^{(h)})$$

$g()$  is sigmoid or ReLU

$$\hat{\mathbf{y}} = f(\mathbf{W}^{(o)\top} \mathbf{h} + \mathbf{w}_0^{(o)})$$

$f()$  is sigmoid/softmax or Linear

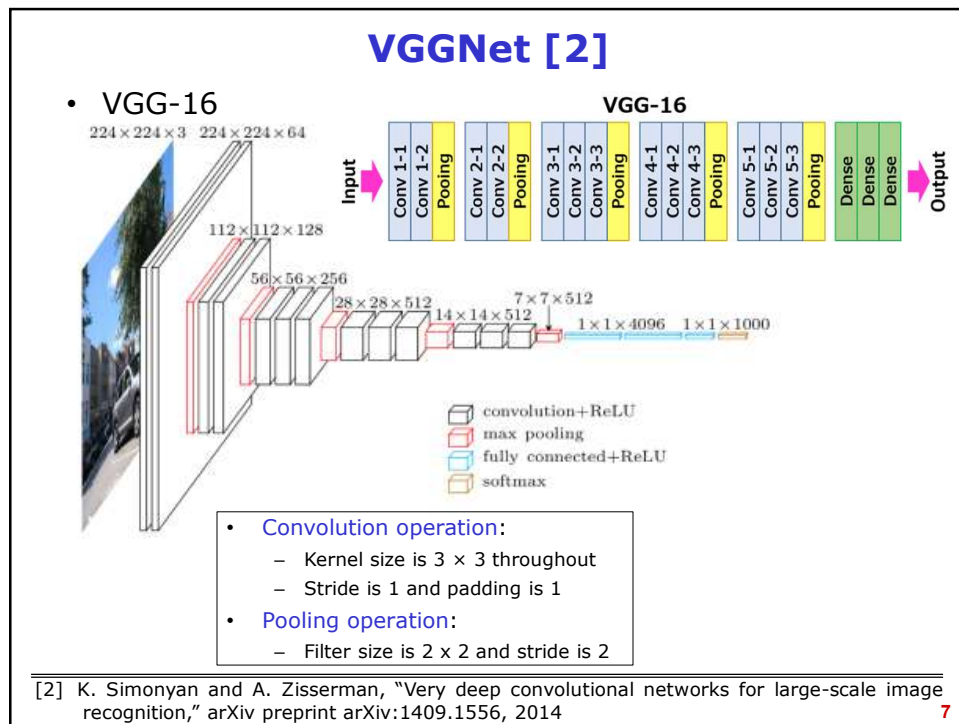
- Weights associated with all the connection between the neurons indicate the parameter of the complex nonlinear discriminant function that the network is trying to approximate
- We train the FCNN using backpropagation
  - Need to compute the gradient of loss function with respect to weight parameters

5

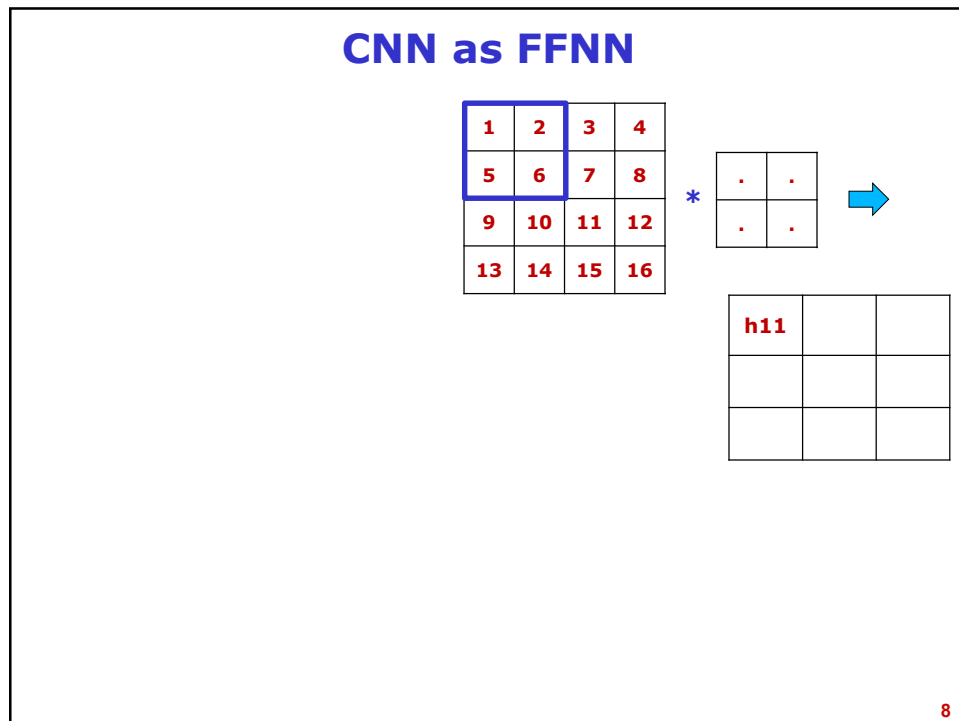
## Convolutional Neural Networks (CNN)

- CNN learn multiple layers of meaningful kernels/filters in addition to learning the weights of the classifier
  - The connections are much sparser
  - Weight sharing
- A CNN can be implemented as a feed-forward neural network
  - Only a few weights are active
  - Rest of the weights are zero
- Each hidden layer is the resultant of convolution operation on the previous layer
  - Rectified linear function (ReLU) is used as activation function on the out put of convolution operation
- It has alternate convolution and pooling layers
- CNN is using backpropagation by considering it as a feed-forward neural network with sparse connection and weight sharing

6

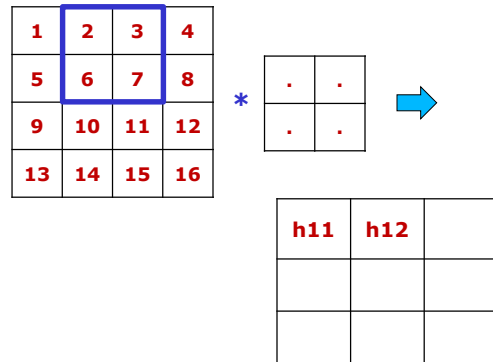


7



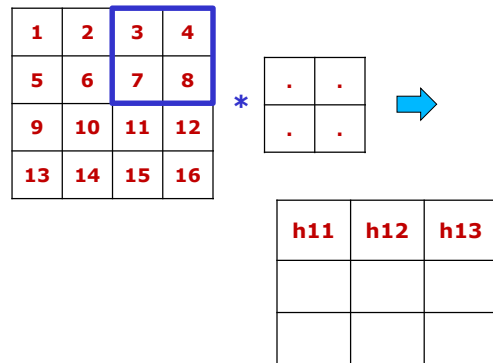
8

## CNN as FFNN



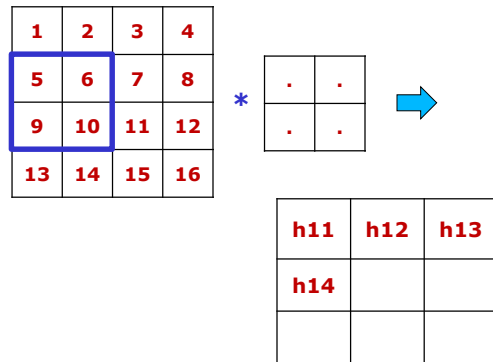
9

## CNN as FFNN



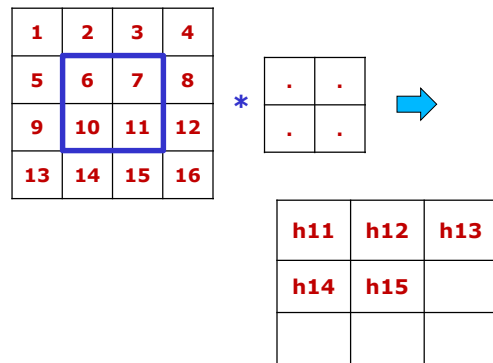
10

## CNN as FFNN



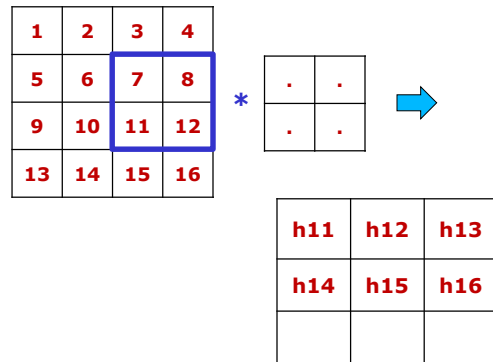
11

## CNN as FFNN



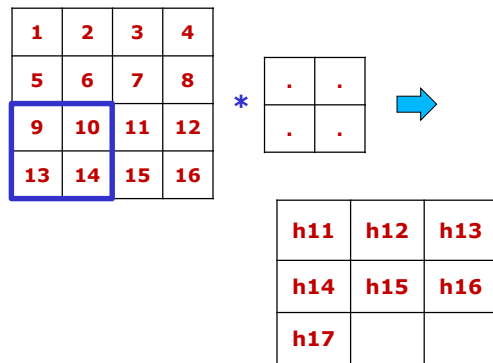
12

## CNN as FFNN



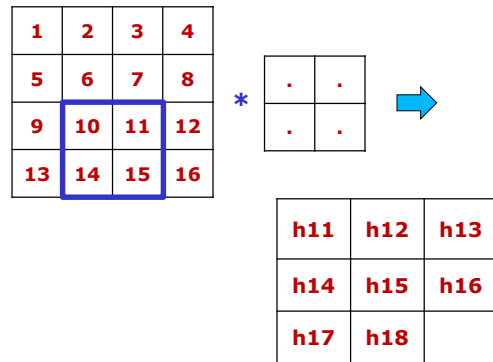
13

## CNN as FFNN



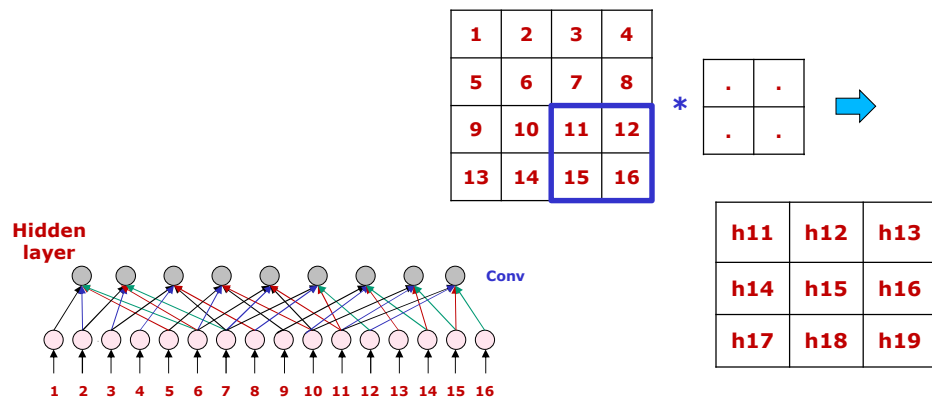
14

## CNN as FFNN



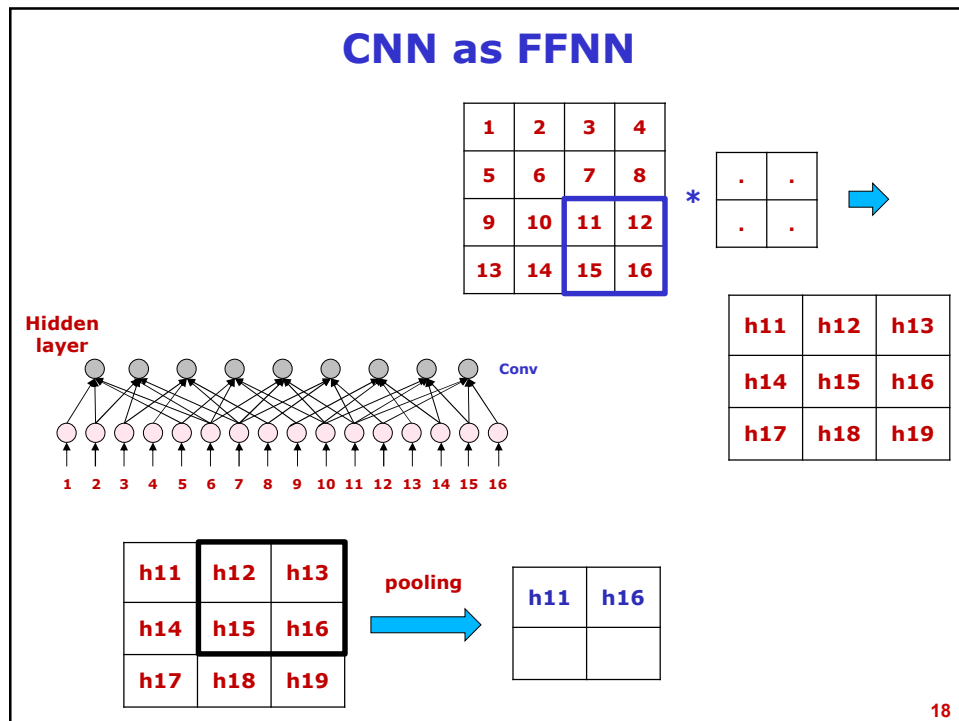
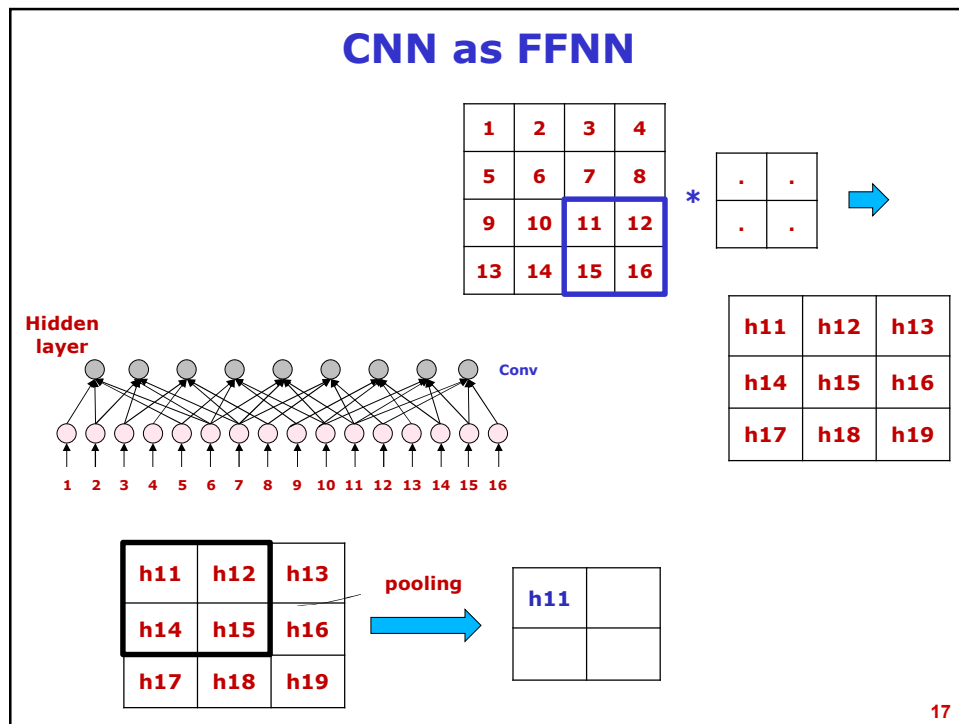
15

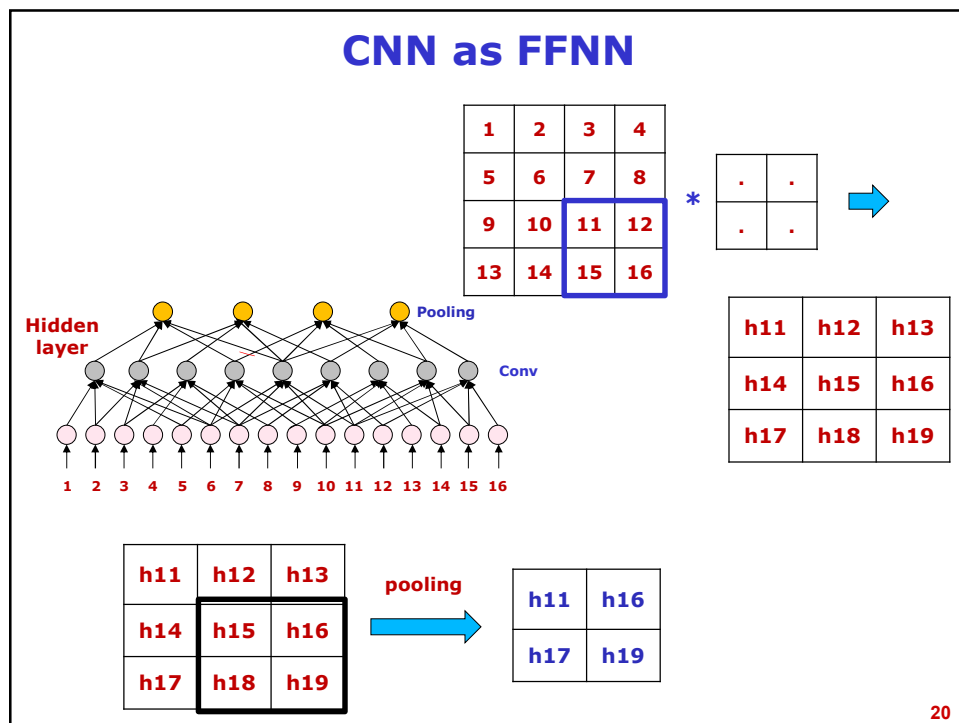
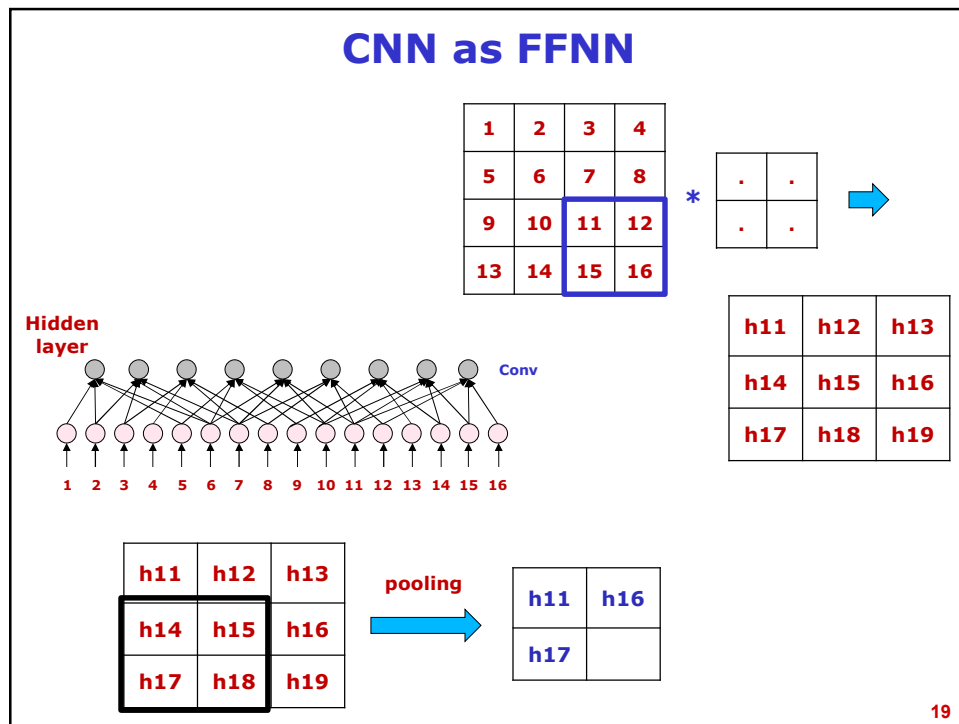
## CNN as FFNN



16







# CNN as FFNN

The diagram illustrates a CNN architecture where the input is a 4x4 grid of values (1-16). This input is processed by a convolutional layer (labeled 'Conv') to produce a 3x3 grid of hidden layer values (h11-h19). The hidden layer values are then processed by a pooling layer (labeled 'Pooling') to produce a 2x2 grid of output values (h11-h19). The output values are then used to classify the input as 'Cat', 'Dog', or 'Tiger'.

The input grid is:

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

The hidden layer values are:

h11	h12	h13
h14	h15	h16
h17	h18	h19

The output values are:

h11	h16
h17	h19

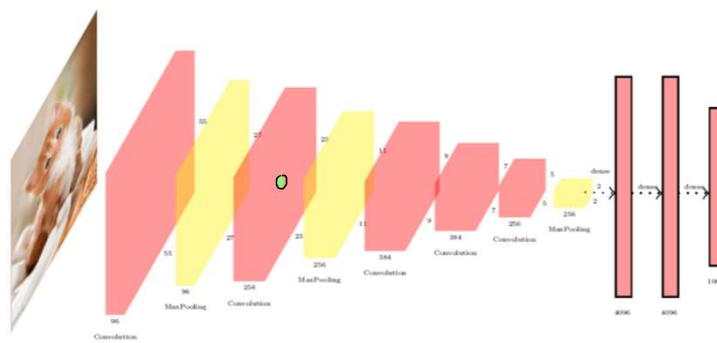
The output values are used to classify the input as 'Cat', 'Dog', or 'Tiger'.

21

## Visualizing Patches which Maximally Activate a Neuron



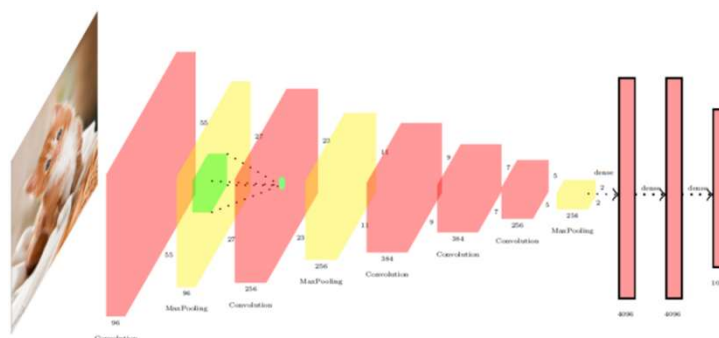
## Visualizing Patches which Maximally Activate a Neuron



- Consider some neurons in a given layer of a CNN
- Feed in images to this CNN and identify the images which cause these neurons to fire

23

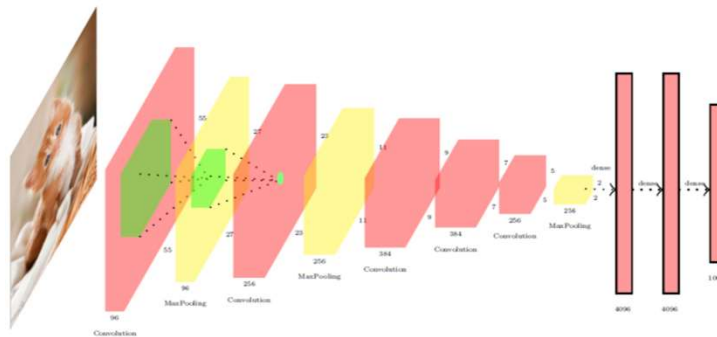
## Visualizing Patches which Maximally Activate a Neuron



- Consider some neurons in a given layer of a CNN
- Feed in images to this CNN and identify the images which cause these neurons to fire
- Then trace back to the patch in the image which causes these neurons to fire

24

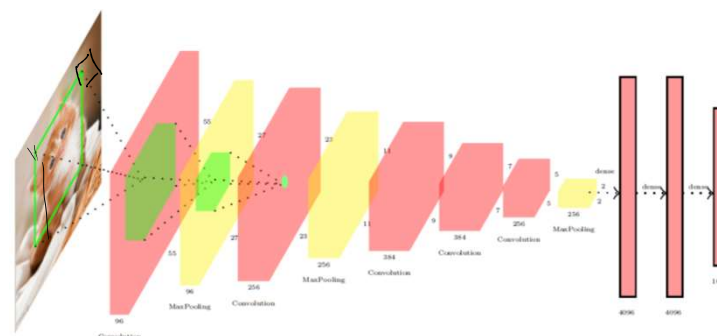
## Visualizing Patches which Maximally Activate a Neuron



- Consider some neurons in a given layer of a CNN
- Feed in images to this CNN and identify the images which cause these neurons to fire
- Then trace back to the patch in the image which causes these neurons to fire

25

## Visualizing Patches which Maximally Activate a Neuron

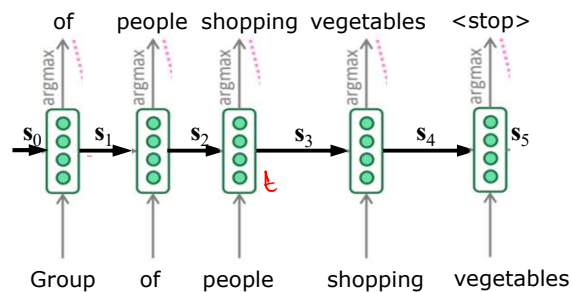


- Consider some neurons in a given layer of a CNN
- Feed in images to this CNN and identify the images which cause these neurons to fire
- Then trace back to the patch in the image which causes these neurons to fire

26

## Recurrent Neural Networks (RNN)

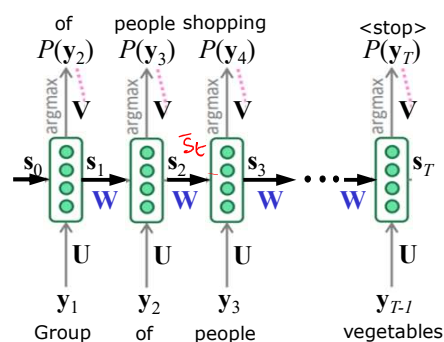
- RNN used for sequential learning problem
  - Each input is dependent on the previous or future input
  - In many applications the input is not of a fixed size
- Consider the problem of language modelling: **Natural sentence generation**
  - Given  $t - i$  words predict the  $t^{\text{th}}$  word
- Example: Generate a sentence – "Group of people shopping vegetables"
- A word **shopping** is predicted given the words **Group, of, people**



27

## Sequence Learning Problem: RNN

- Sequence learning:** More formally, given  $y_1, y_2, \dots, y_{t-1}$  we want to find  $\hat{y} = \arg \max P(y_t = j | y_1, y_2, \dots, y_{t-1})$ 
  - where  $j \in \mathcal{V}$  and  $\mathcal{V}$  is the set of all the words in vocabulary
- Let us denote  $P(y_t = j | y_1, y_2, \dots, y_{t-1})$  as  $P(y_t = j | (y)_1^{t-1})$

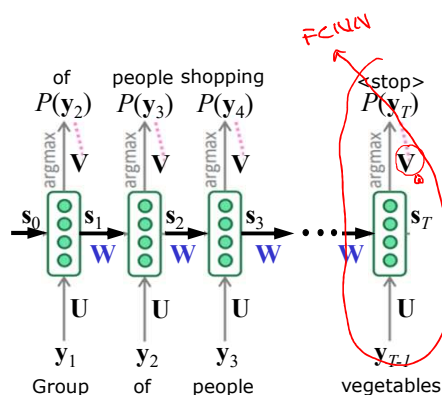


- Using RNN:
  - $P(y_t = j | (y)_1^{t-1}) = \text{softmax}(\mathbf{V}\mathbf{s}_t + c)_j$
  - $\mathbf{s}_t$  is the **hidden representation** at time step  $t$
  - Recurrent connections ensure that information about sequence  $y_1, y_2, \dots, y_{t-1}$  is embedded in  $\mathbf{s}_t$
  - Hence,
    - $P(y_t = j | (y)_1^{t-1}) = P(y_t = j | \mathbf{s}_t)$

28

## Sequence Learning Problem: RNN

- **Sequence learning:** More formally, given  $y_1, y_2, \dots, y_{t-1}$  we want to find  $\hat{y} = \arg \max P(y_t = j | y_1, y_2, \dots, y_{t-1})$ 
  - where  $j \in \mathcal{V}$  and  $\mathcal{V}$  is the set of all the words in vocabulary
- Let us denote  $P(y_t = j | y_1, y_2, \dots, y_{t-1})$  as  $P(y_t = j | (y)_1^{t-1})$



- Using RNN:
 
$$P(y_t = j | s_t) = \text{softmax}(\mathbf{V}\mathbf{s}_t + c)_j$$
- Recurrent connections ensure that information about sequence  $y_1, y_2, \dots, y_{t-1}$  is embedded in  $s_t$ 

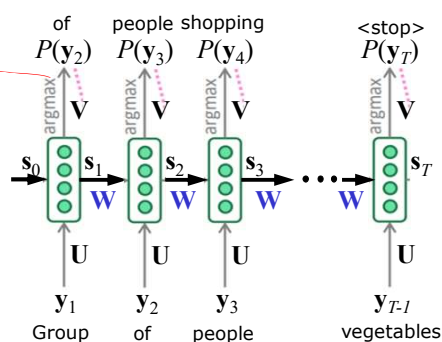
$$s_t = \text{sigmoid}(\mathbf{U}\mathbf{y}_t + \mathbf{W}\mathbf{s}_{t-1} + b)$$

$$s_t = \text{RNN}(s_{t-1}, y_t)$$

29

## Language Modelling Problem: Natural Sentence Generation : RNN

- **Data:** All sentences from any large corpus (say Wikipedia)
- Each word in the vocabulary is represented as  $d$ -dimensional word vector (example: word-to-vec)
- RNN is trained using **backpropagation through time** (BPTT)

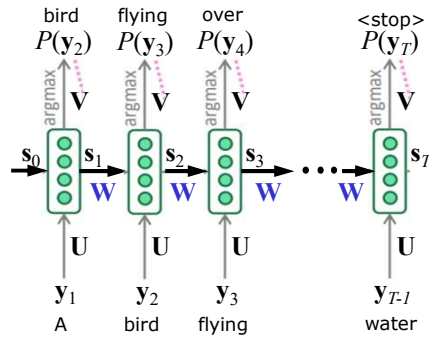


- One can also use LSTM or GRU in the place of RNN

30

## Neural Image Caption Generation

- So far we have seen how to generate a sentence given previous words
- Now, we want to generate a sentence given an image

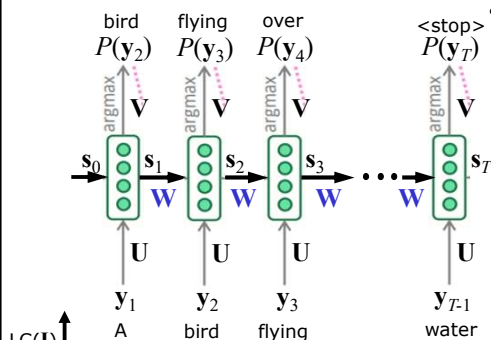


A bird flying over a body of water

- We are now interested in  $P(y_t = j | y_1, y_2, \dots, y_{t-1}, \mathbf{I})$  instead of  $P(y_t = j | y_1, y_2, \dots, y_{t-1})$ 
  - where  $\mathbf{I}$  is an image
- Usually information in the image is encoded in a feature vector

31

## Neural Image Caption Generation

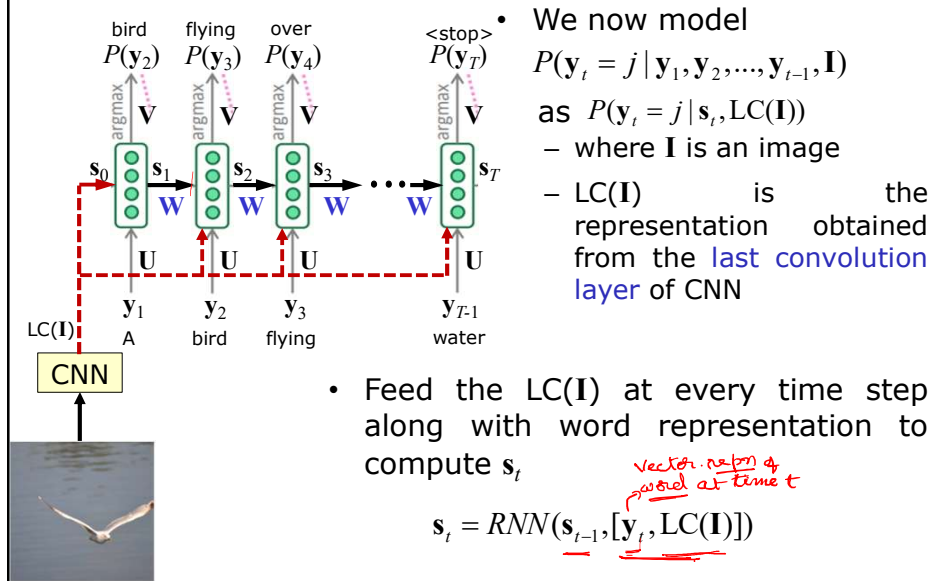


- We now model  $P(y_t = j | y_1, y_2, \dots, y_{t-1}, \mathbf{I})$  as  $P(y_t = j | s_t, \text{LC}(\mathbf{I}))$ 
  - where  $\mathbf{I}$  is an image
  - $\text{LC}(\mathbf{I})$  is the representation obtained from the last convolution layer of CNN
- Feed the  $\text{LC}(\mathbf{I})$  at every time step along with word representation to compute  $s_t$

32



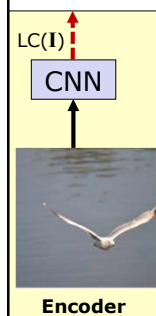
## Neural Image Caption Generation



33

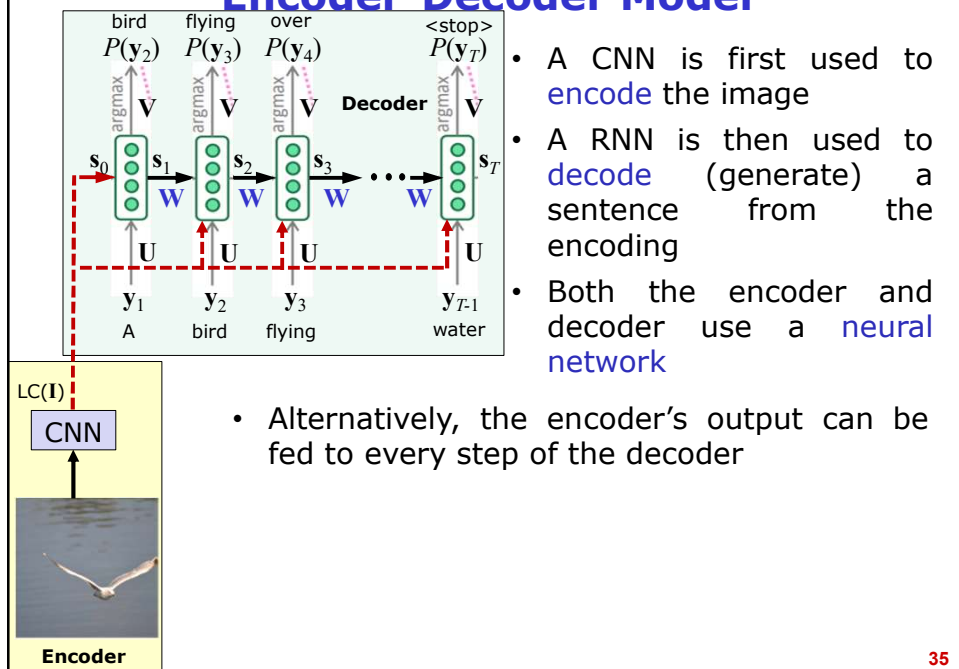
## Encoder-Decoder Model

- A CNN is first used to **encode** the image



34

## Encoder-Decoder Model



35

## More Applications of Encoder-Decoder Models

- **Machine Translation:**
  - Translating sentence in one language to another
  - Encoder: RNN
  - Decoder: RNN
- **Transliteration:**
  - Translating the script of one language to script of another language
  - Encoder: RNN
  - Decoder: RNN
- **Image Question Answering:**
  - Given the image and a question (sentence), generate answer (word)
  - Encoder: CNN + RNN
  - Decoder: FCNN

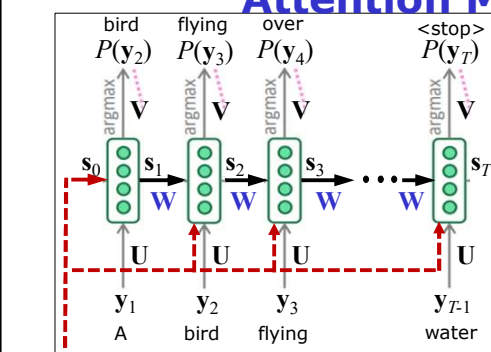
36

## More Applications of Encoder-Decoder Models

- Document Summarization:
  - Generating a summary of a document
  - Encoder: RNN
  - Decoder: RNN
- Video Captioning:
  - Generate sentence given video
  - Encoder: CNN-RNN
  - Decoder: RNN
- And many more ...

37

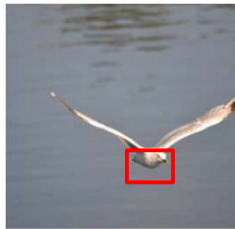
## Attention Mechanism



- Encoder-decoder models can be made even more expressive by adding an “attention” mechanism
- Let us motivate the task of attention with the help of **image captioning**
- The encoder reads the image only once and encodes it
  - Embedding from last convolution layer of CNN
- At each timestep the decoder uses this embedding to produce a new word

38

## Attention Mechanism: Image Captioning

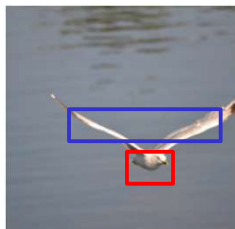


- Humans try to produce each word in the output by **focusing only on certain objects (concepts)** in the input image
- Example:

A **bird** flying over a body of water

39

## Attention Mechanism: Image Captioning

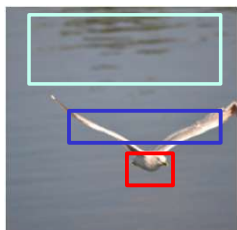


- Humans try to produce each word in the output by **focusing only on certain objects (concepts)** in the input image
- Example:

A **bird** **flying** over a body of water

40

## Attention Mechanism: Image Captioning



- Humans try to produce each word in the output by **focusing only on certain objects (concepts)** in the input image
- Example:
- Essentially at each time step we come up with **a distribution (weights)** on the input concepts (objects)

A bird flying over a body of water

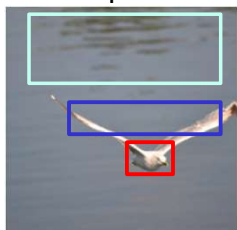
- This distribution tells us **how much attention to pay to each objects** in input at each time step
- Ideally, at each time step we should **feed only this relevant information** (i.e. encodings of relevant objects) to the decoder

• $t_1$ : A	[ 1 0 0 ]
• $t_2$ : bird	[ 0.5 0.5 0 ]
• $t_3$ : flying	[ 0 1 0 ]
• $t_4$ : over	[ 0 1 0 ]
• $t_5$ : a	[ 0 0 1 ]
• $t_6$ : body	[ 0 0 1 ]
• $t_7$ : of	[ 0 0 1 ]
• $t_8$ : water	[ 0 0 1 ]

41

## Attention Mechanism: Image Captioning

- Humans try to produce each word in the output by **focusing only on certain objects (concepts)** in the input image
- Example:



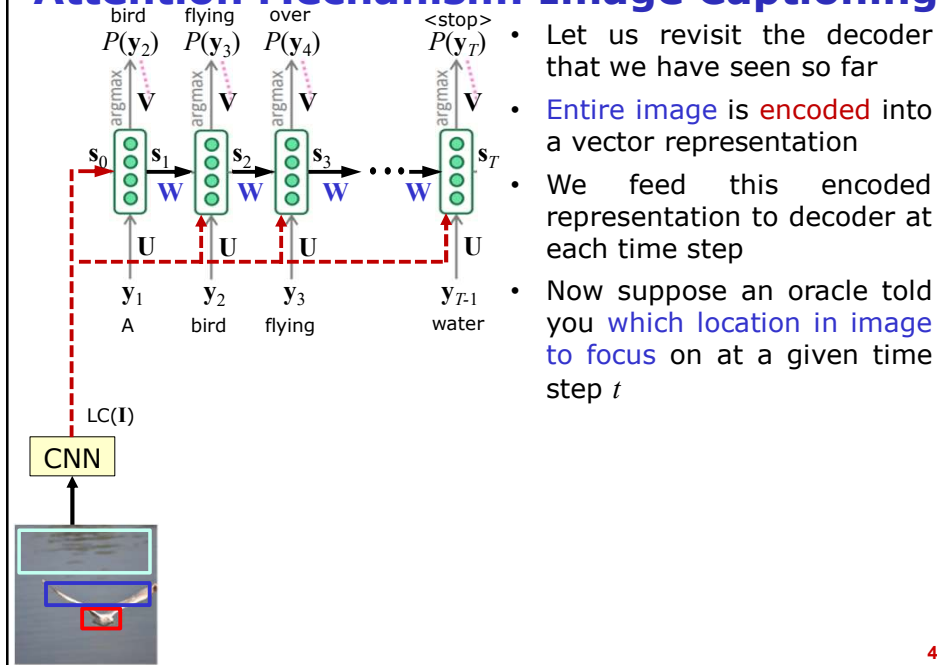
A bird flying over a body of water



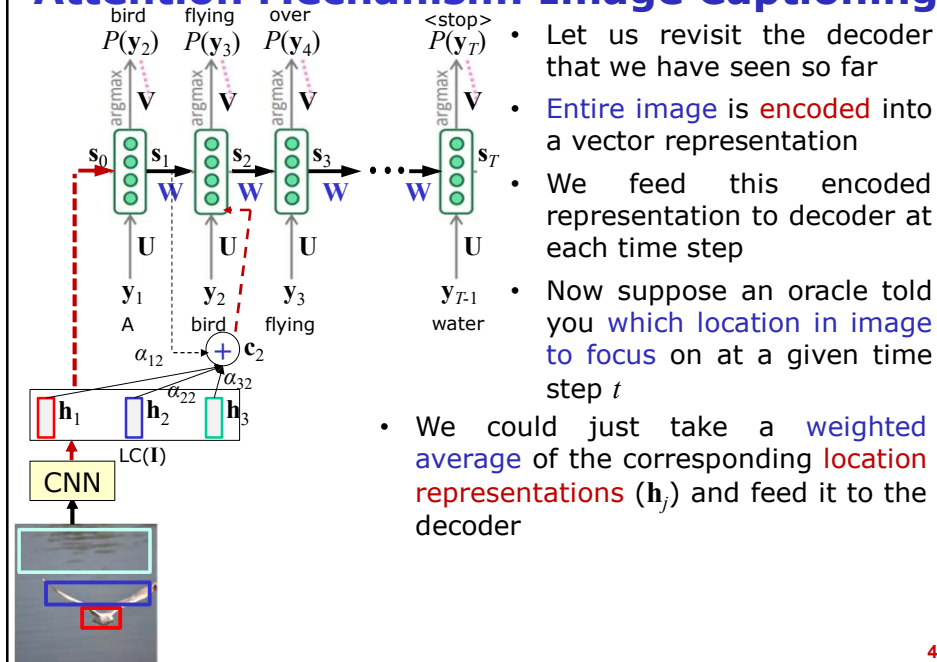
A group of people sitting on a boat in the water

42

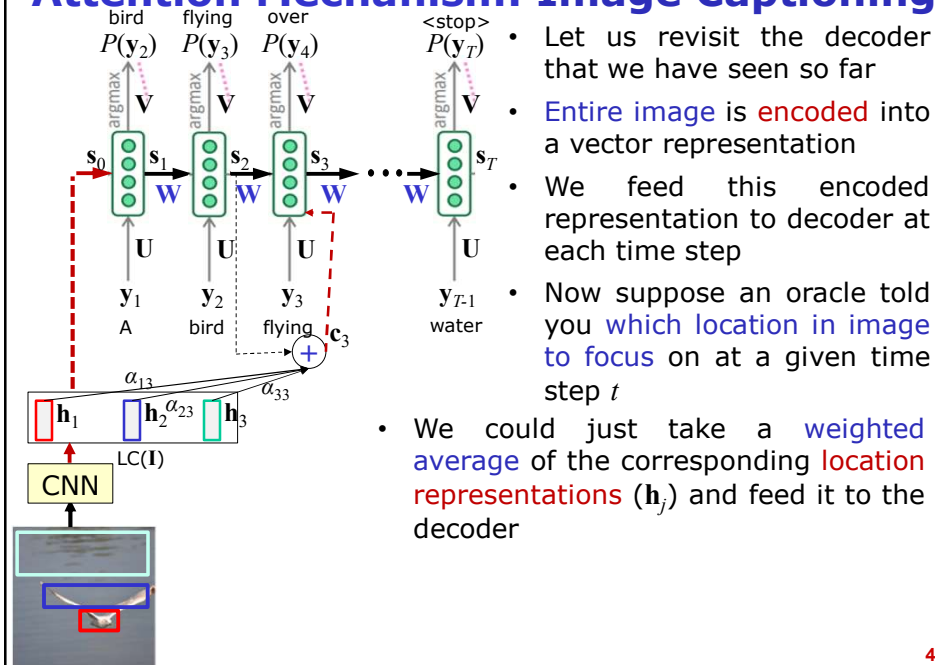
## Attention Mechanism: Image Captioning



## Attention Mechanism: Image Captioning

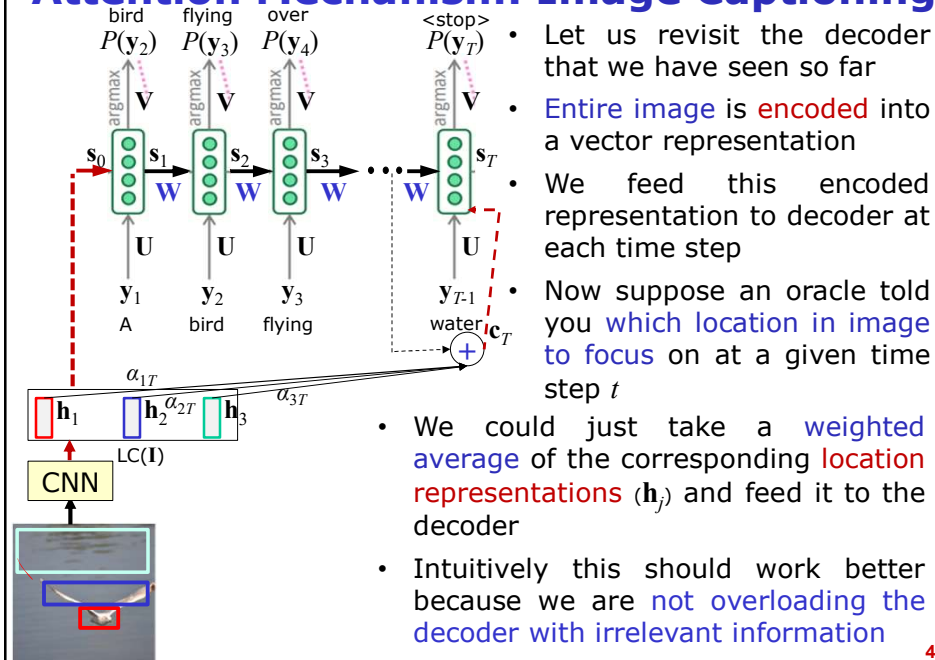


## Attention Mechanism: Image Captioning



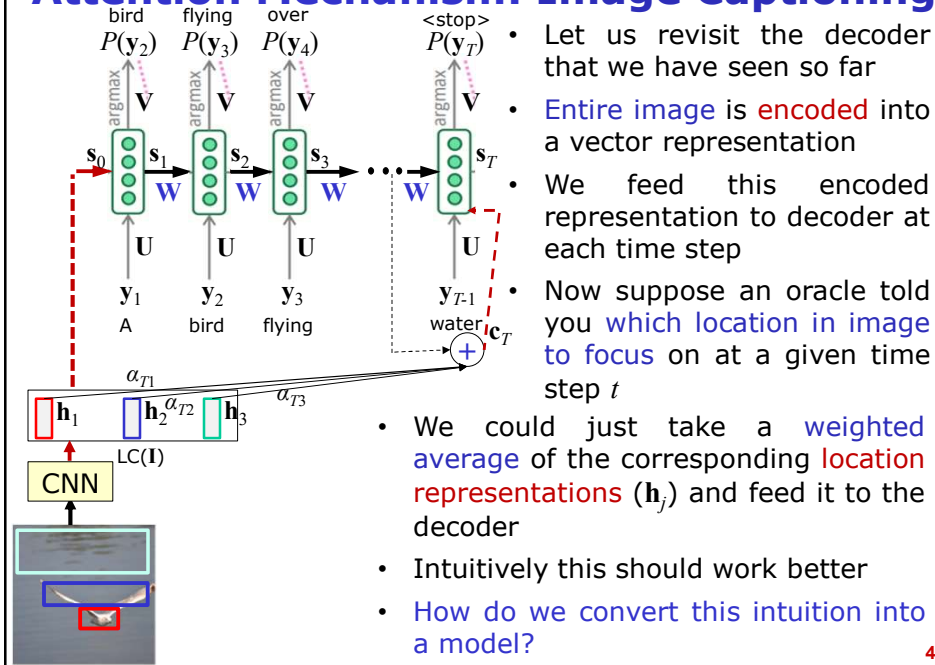
45

## Attention Mechanism: Image Captioning



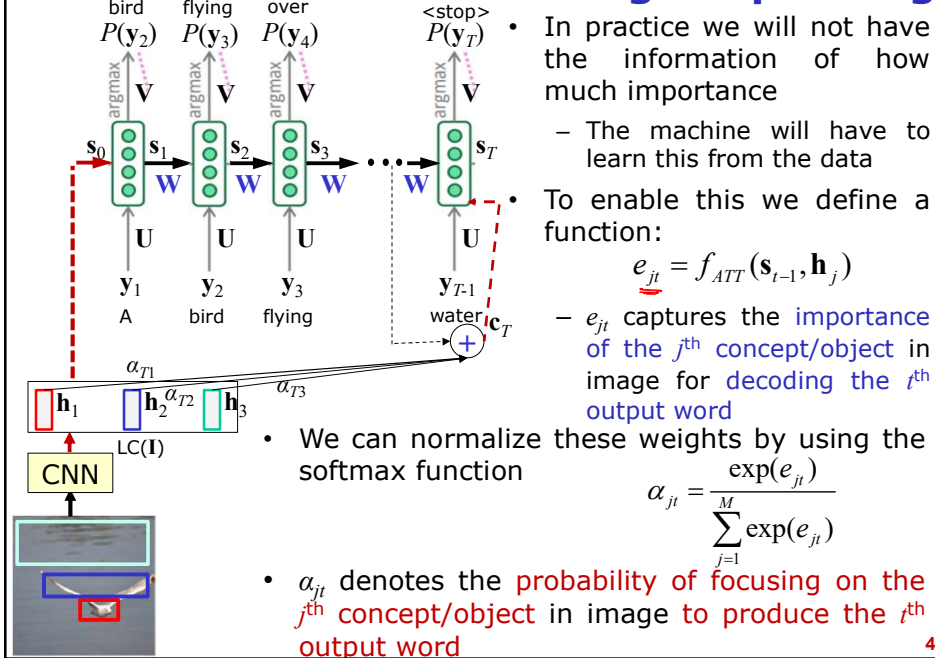
46

## Attention Mechanism: Image Captioning



47

## Attention Mechanism: Image Captioning



48



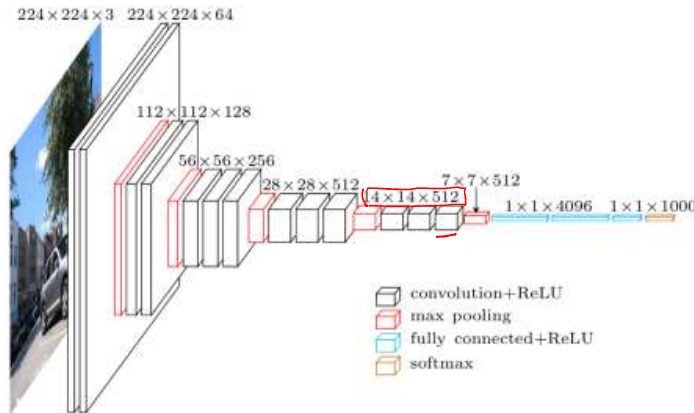


## Attention Mechanism: Image Captioning



## Learning Attention Over Image Location

- Consider VGG16 network to encode image.

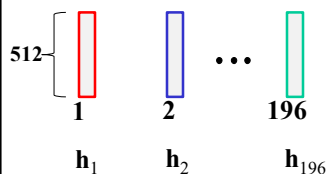
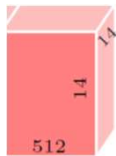


- Output of last convolution layer is a  $14 \times 14 \times 512$  feature map
- We could think of this as 196 (i.e.  $14 \times 14$ ) locations (each having a 512 dimensional representation)

51

## Learning Attention Over Image Location

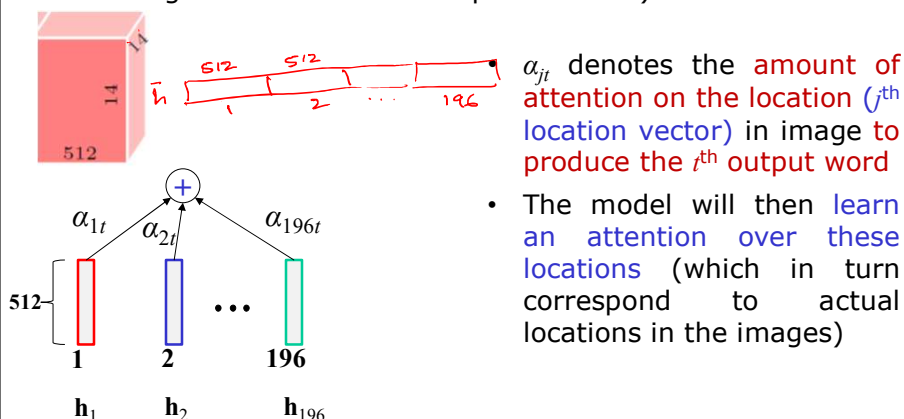
- Consider VGG16 network to encode image.
- Output of last convolution layer is a  $14 \times 14 \times 512$  feature map
- We could think of this as 196 (i.e.  $14 \times 14$ ) locations (each having a 512 dimensional representation)



52

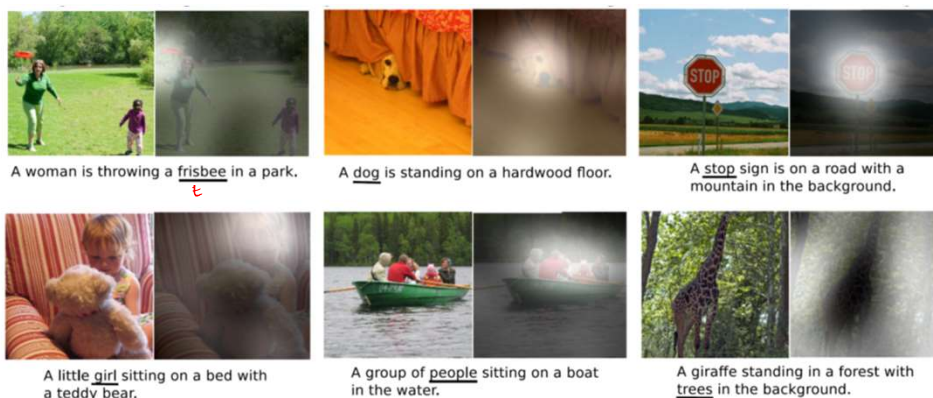
## Learning Attention Over Image Location

- Consider VGG16 network to encode image.
- Output of last convolution layer is a 14x14x512 feature map
- We could think of this as 196 (i.e. 14x14) locations (each having a 512 dimensional representation)



53

## Illustrations: Attention Over Images

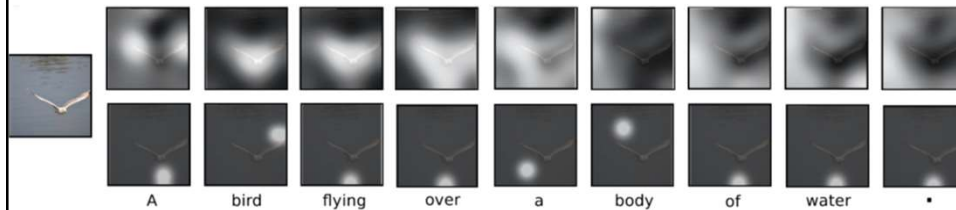


Examples of the attention-based model attending to the correct object (white indicates the attended regions, underlines indicates the corresponding word) [3]

- [3] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, Yoshua Bengio, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", in *Proceedings of the 32nd International Conference on Machine Learning*, PMLR vol. 37, pp. 2048-2057, 2015.

54

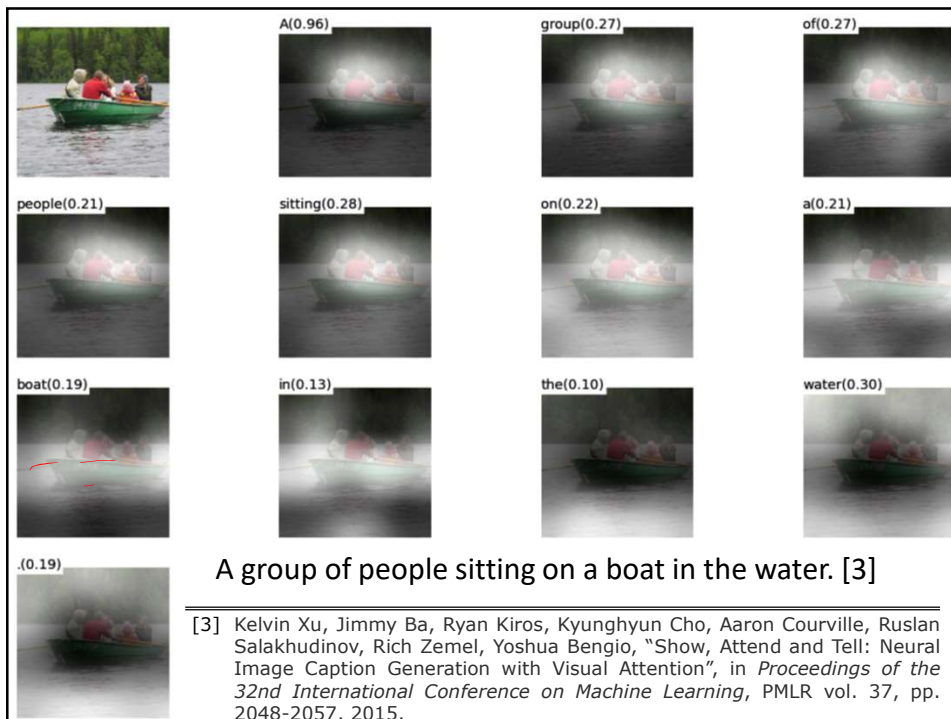
## Illustrations: Attention Over Images



**Attention over time.** As the model generates each word, its attention changes to reflect the relevant parts of the image. (white indicates the attended regions) [3]

- [3] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, Yoshua Bengio, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", in *Proceedings of the 32nd International Conference on Machine Learning*, PMLR vol. 37, pp. 2048-2057, 2015.

55



- [3] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, Yoshua Bengio, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", in *Proceedings of the 32nd International Conference on Machine Learning*, PMLR vol. 37, pp. 2048-2057, 2015.

## Summary

- Attention mechanism in encoder-decoder models
- Encoder-decoder model
  - Encoder first used to **encode** the input
  - A decoder is then used to **decode** (generate) a output from the encoding
- Encoder-decoder models can be made even more expressive by adding an “**attention**” mechanism
- A model will then **learn an attention** over input to generate output
  - Attention is seen as probability of portion of input responsible for generating output

57

## Text Books

1. Ian Goodfellow, Yoshua Bengio and Aaron Courville, **Deep learning**, MIT Press, Available online: <http://www.deeplearningbook.org>, 2016
2. Charu C. Aggarwal, **Neural Networks and Deep Learning**, Springer, 2018
3. B. Yegnanarayana, **Artificial Neural Networks**, Prentice-Hall of India, 1999.
4. Satish Kumar, **Neural Networks - A Class Room Approach**, Second Edition, Tata McGraw-Hill, 2013.
5. S. Haykin, **Neural Networks and Learning Machines**, Prentice Hall of India, 2010.
6. C. M. Bishop, **Pattern Recognition and Machine Learning**, Springer, 2006.
7. J. Han and M. Kamber, **Data Mining: Concepts and Techniques**, Third Edition, Morgan Kaufmann Publishers, 2011.
8. S. Theodoridis and K. Koutroumbas, **Pattern Recognition**, Academic Press, 2009.

58