

Deep Learning for Computer Vision

Vision and Language: Image Captioning

Vineeth N Balasubramanian

Department of Computer Science and Engineering
Indian Institute of Technology, Hyderabad

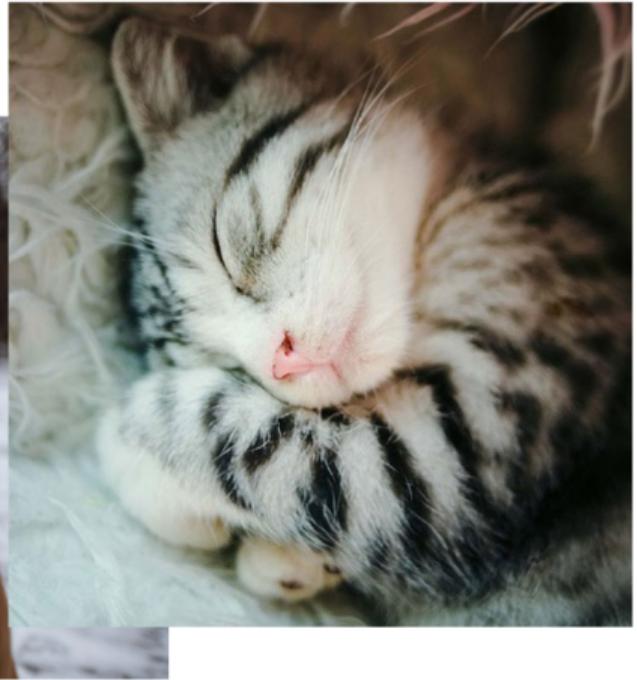
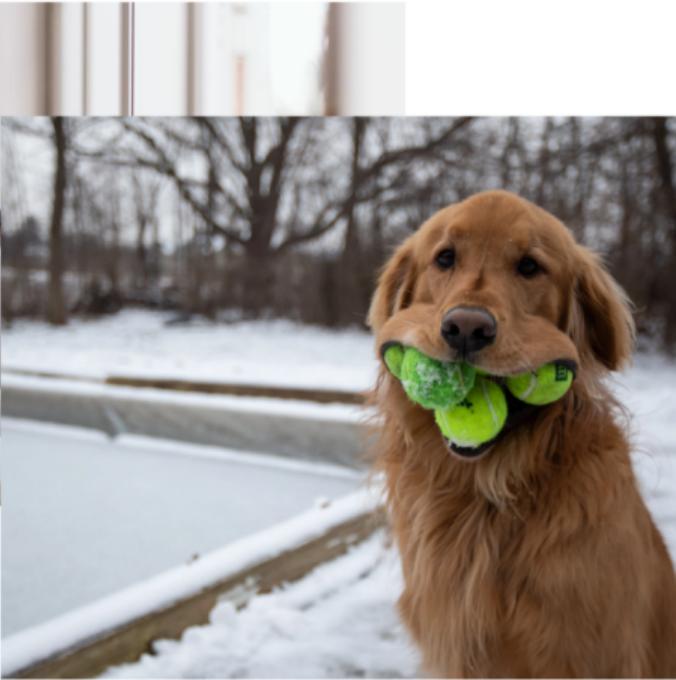


Review: Are autoencoders (AE) and PCA connected?

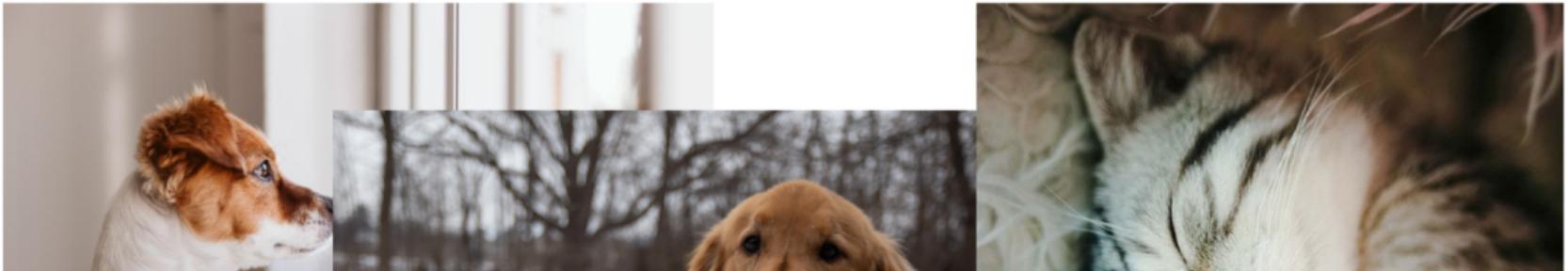
Review: Are autoencoders (AE) and PCA connected?

- Yes!

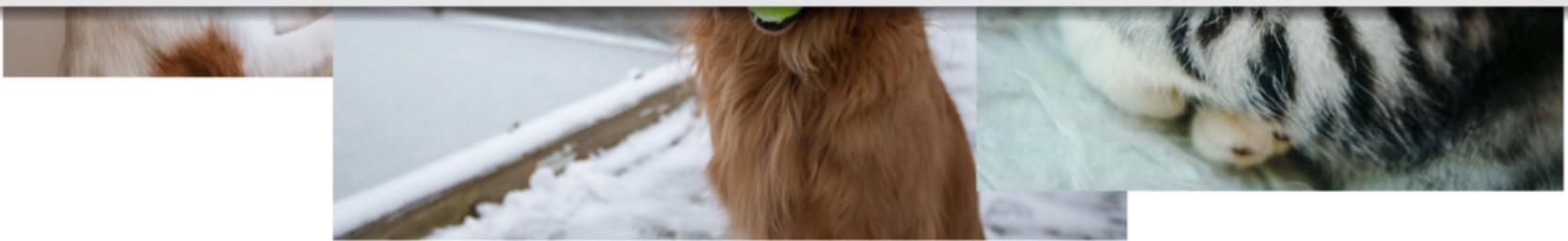
Describe these images



Describe these images



- How can we understand what is happening just by looking at a single image ?
- Can we make a computer do the same?



How to make a computer describe an image?

Some Method

How to make a computer describe an image?



Some Method



How to make a computer describe an image?



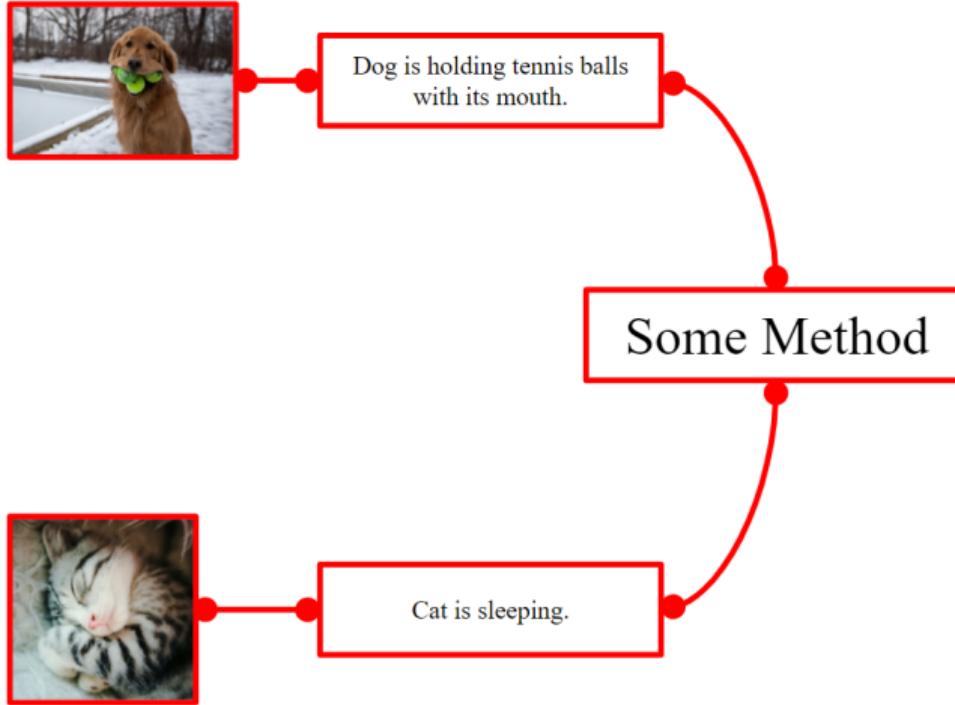
Dog is holding tennis balls
with its mouth.

Some Method

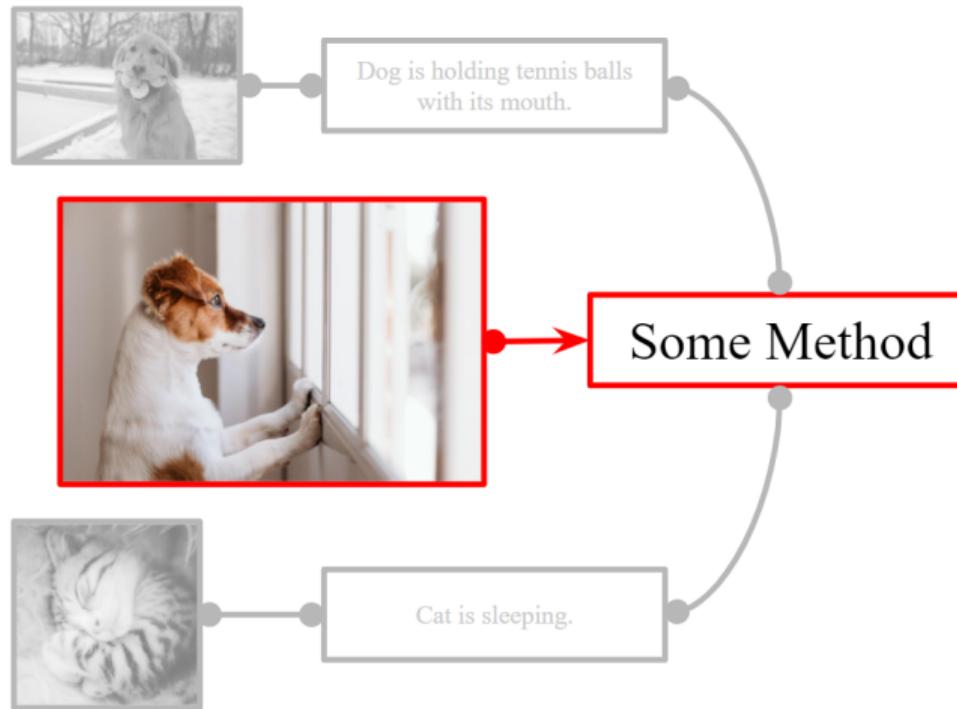


Cat is sleeping.

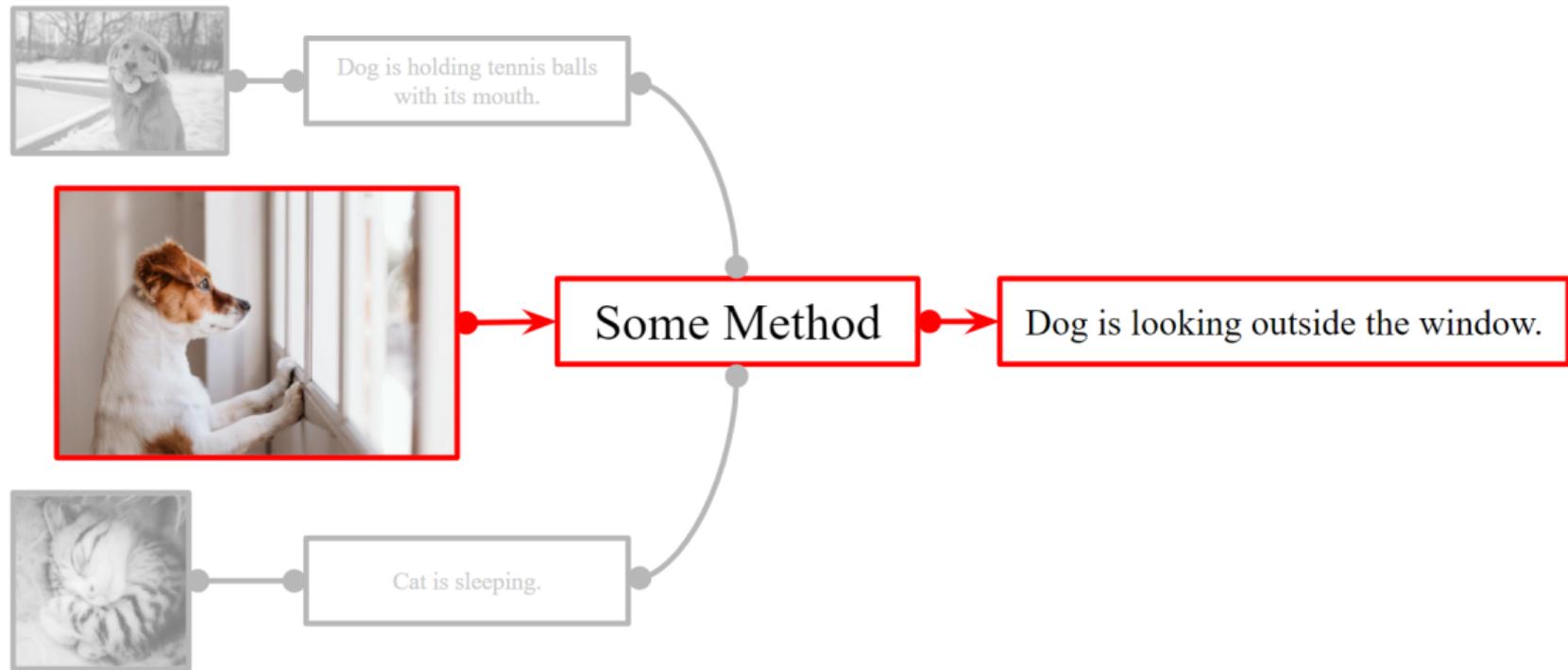
How to make a computer describe an image?



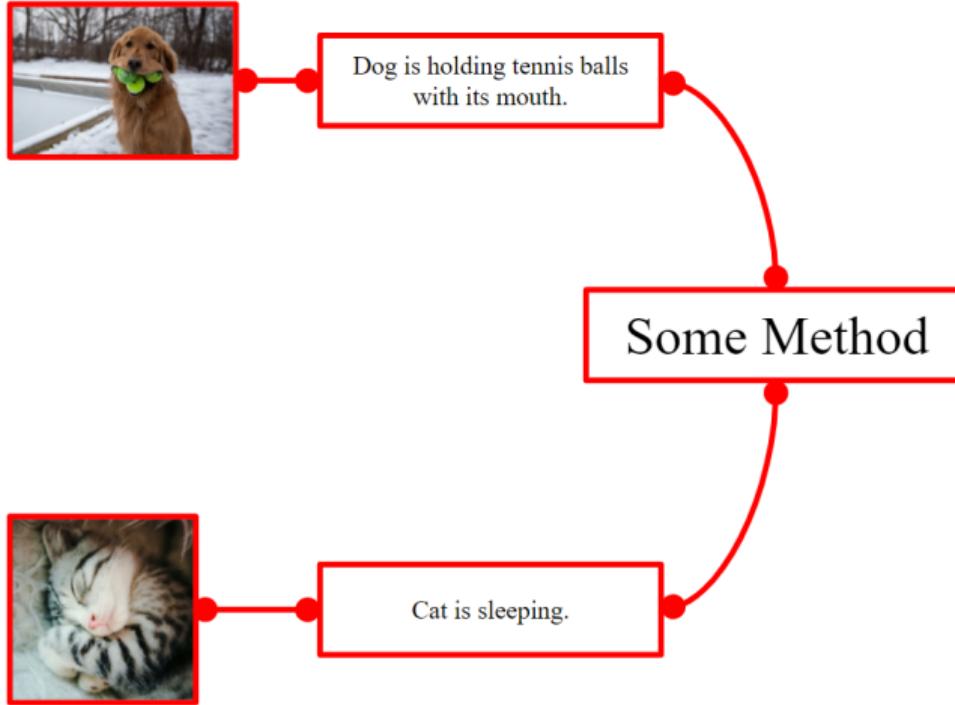
How to make a computer describe an image?



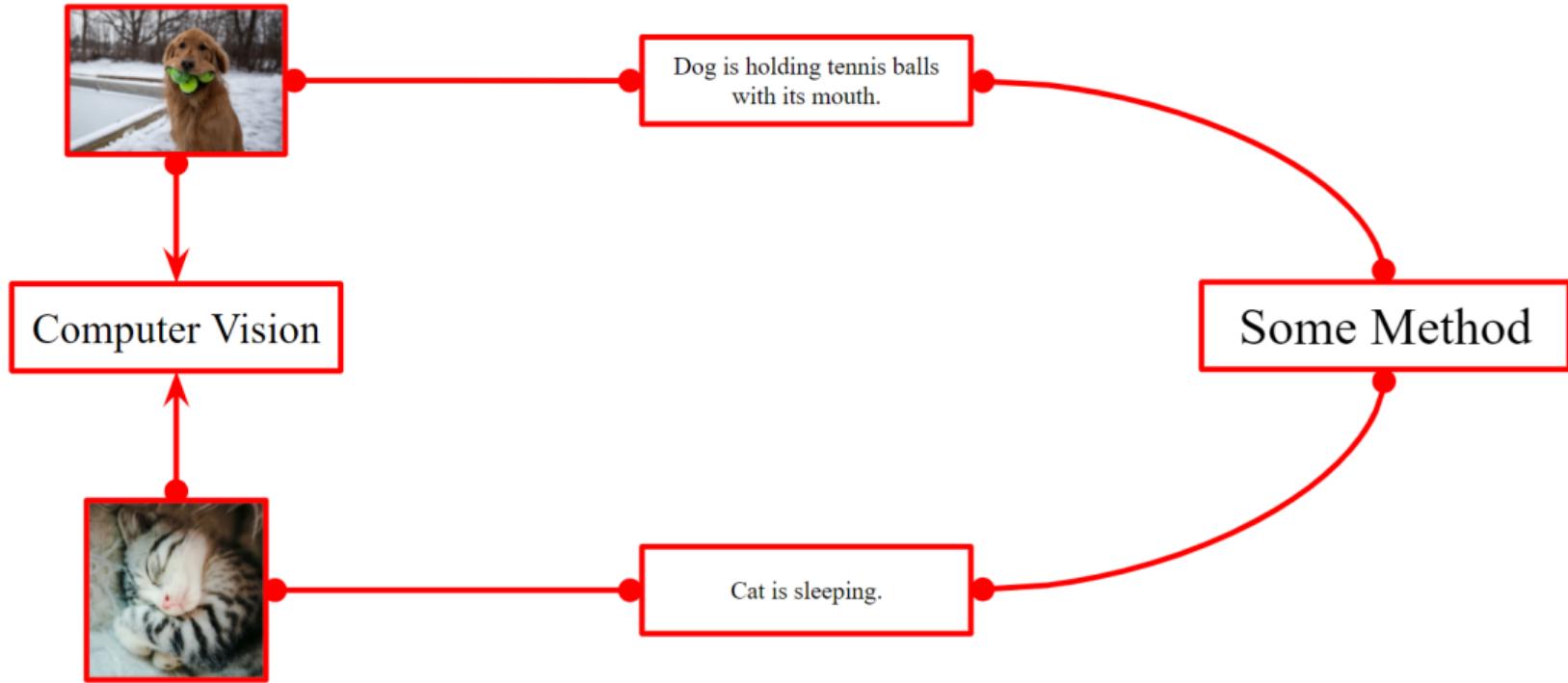
How to make a computer describe an image?



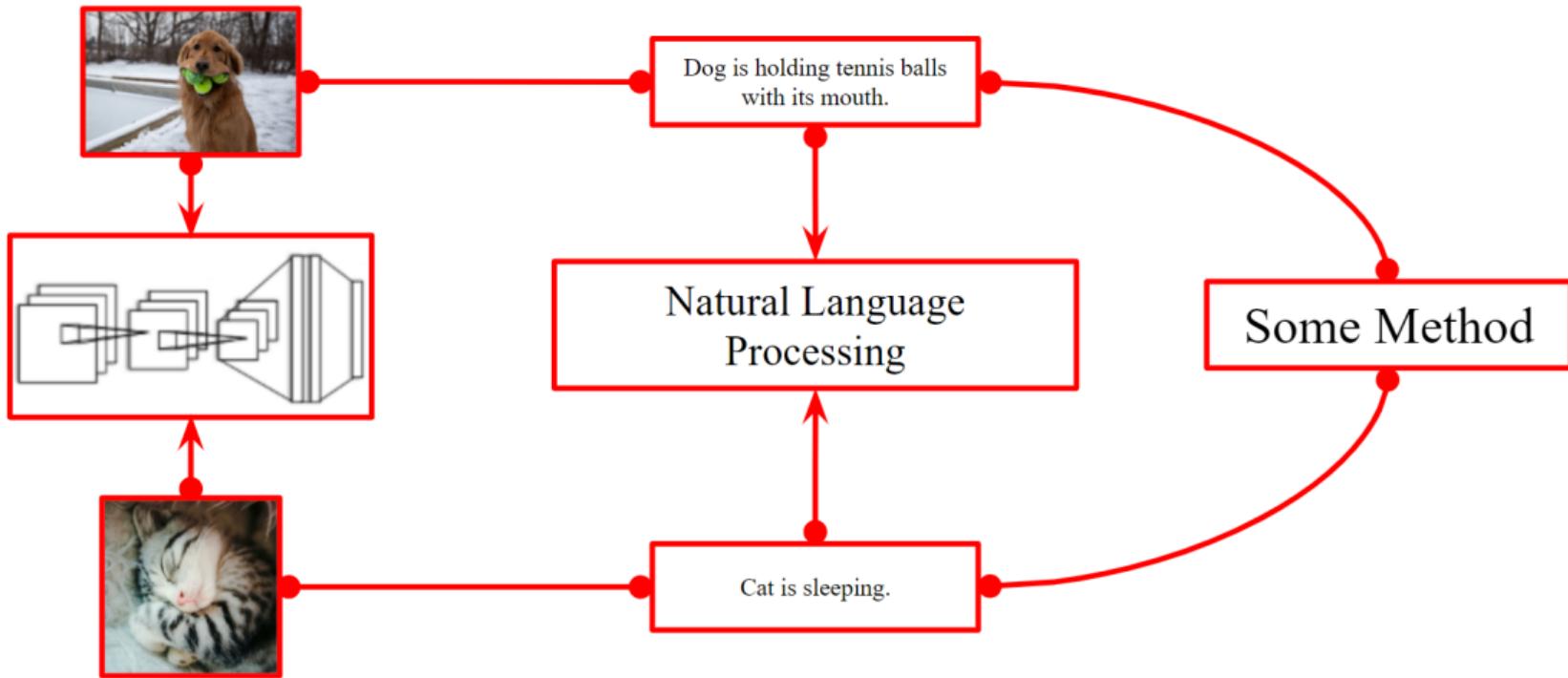
How to make a computer describe an image?



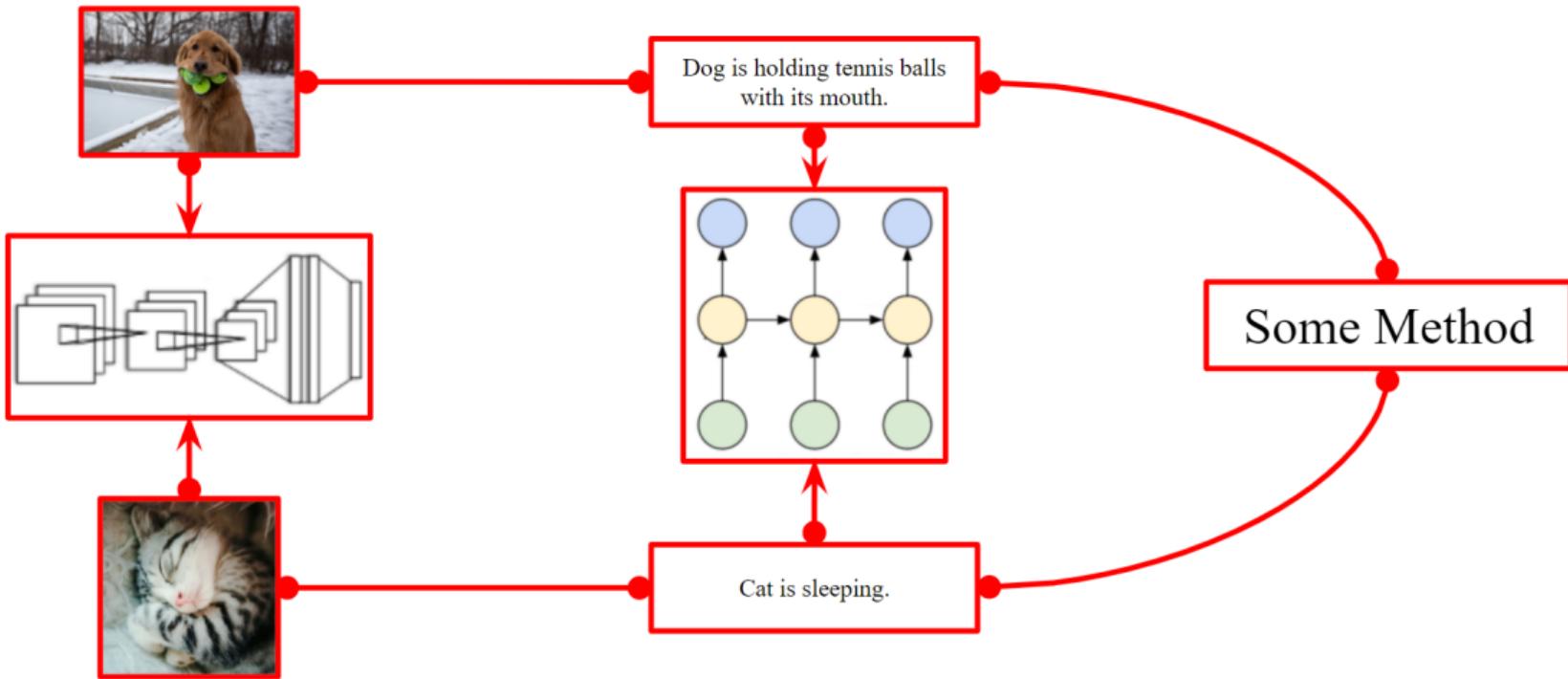
How to make a computer describe an image?



How to make a computer describe an image?



How to make a computer describe an image?



How to make a computer describe an image?

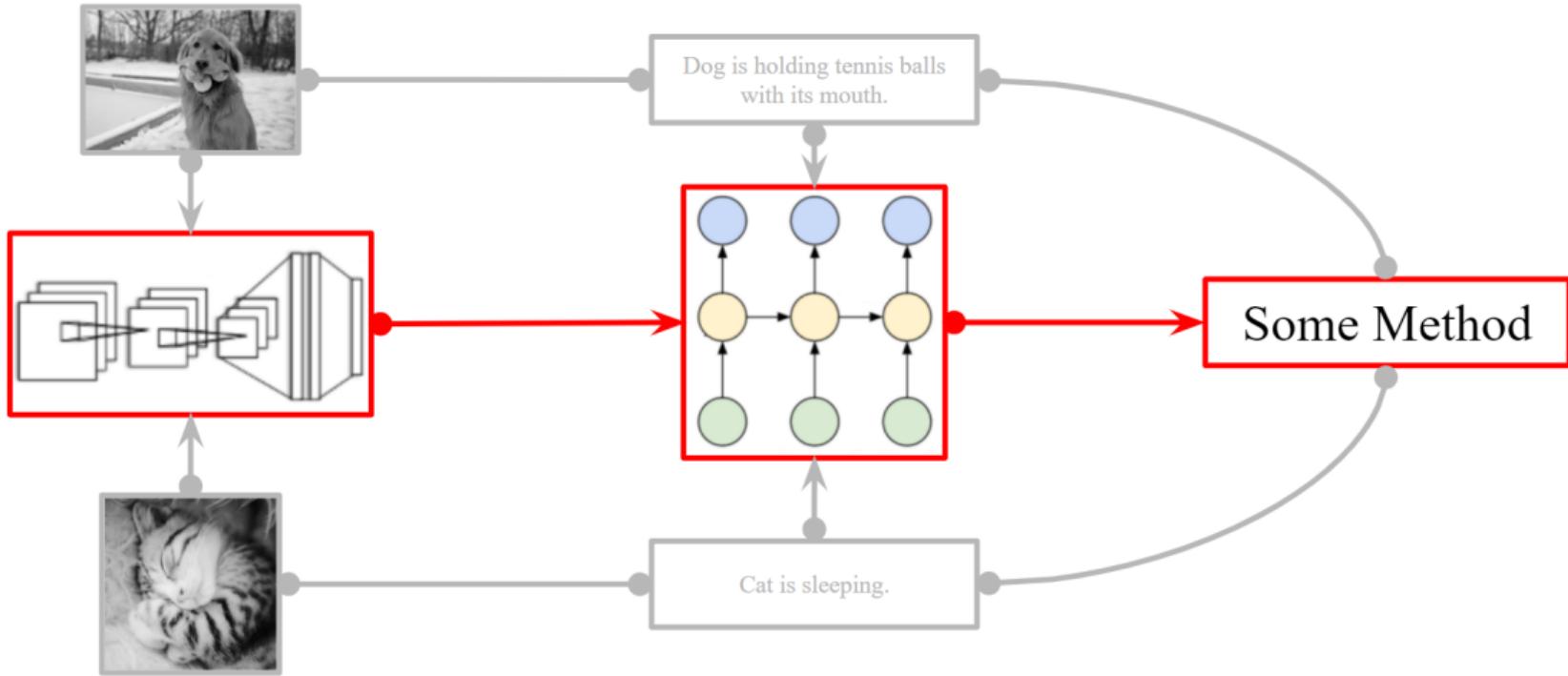


Image Captioning: Training

image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

FC-4096

FC-1000

softmax

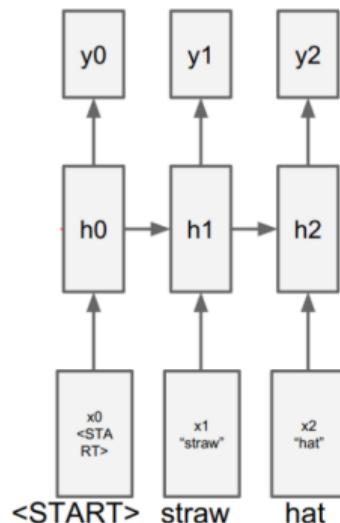


“straw hat”

Credit: Karpathy et al, Deep visual-semantic alignments for generating image descriptions, CVPR 2015

Image Captioning: Training

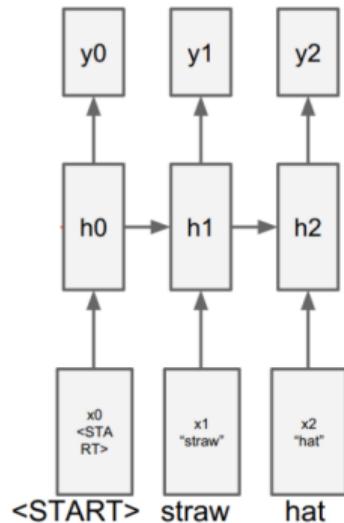
image
conv-64
conv-64
maxpool
conv-128
conv-128
maxpool
conv-256
conv-256
maxpool
conv-512
conv-512
maxpool
conv-512
conv-512
maxpool
FC-4096
FC-4096
FC-1000
softmax



“straw hat”

Credit: Karpathy et al, Deep visual-semantic alignments for generating image descriptions, CVPR 2015

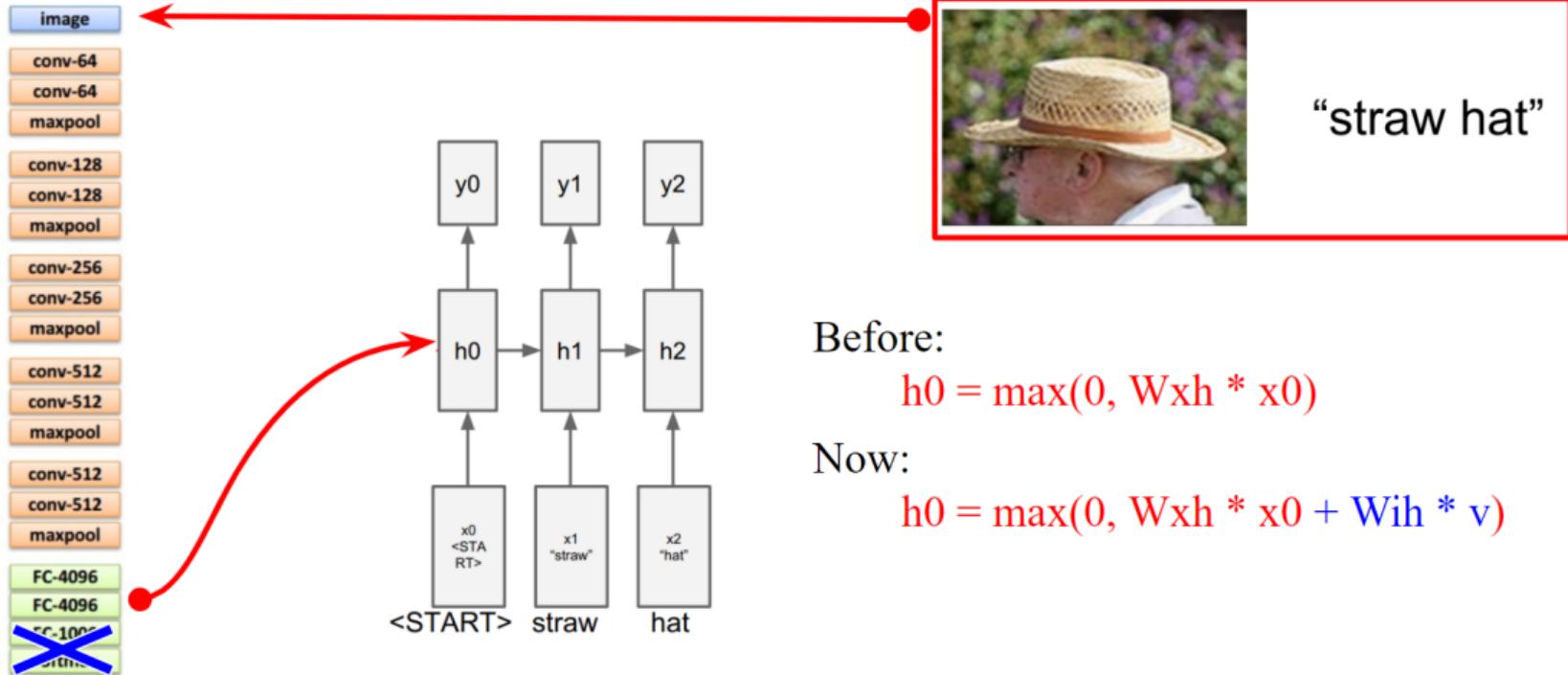
Image Captioning: Training



"straw hat"

Credit: Karpathy et al, Deep visual-semantic alignments for generating image descriptions, CVPR 2015

Image Captioning: Training



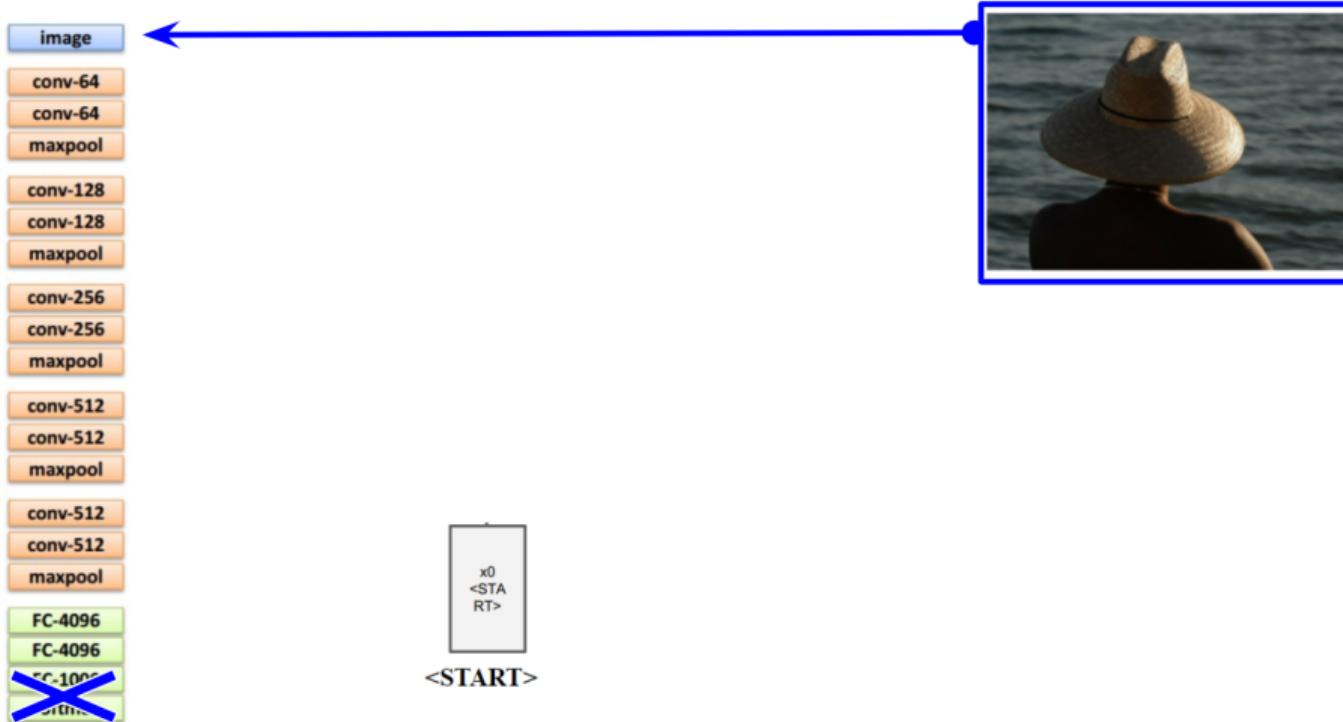
Credit: Karpathy et al, Deep visual-semantic alignments for generating image descriptions, CVPR 2015

Image Captioning: Inference (Test Time)



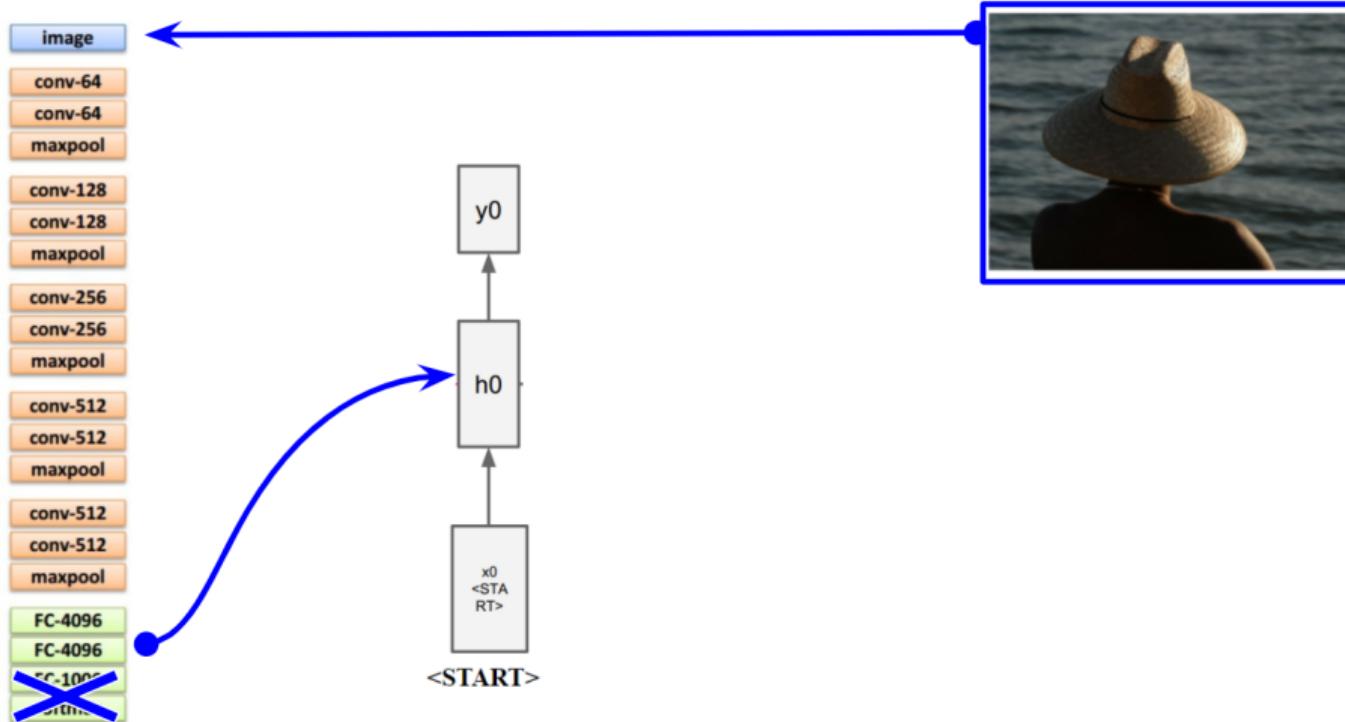
Credit: Karpathy et al, Deep visual-semantic alignments for generating image descriptions, CVPR 2015

Image Captioning: Inference (Test Time)



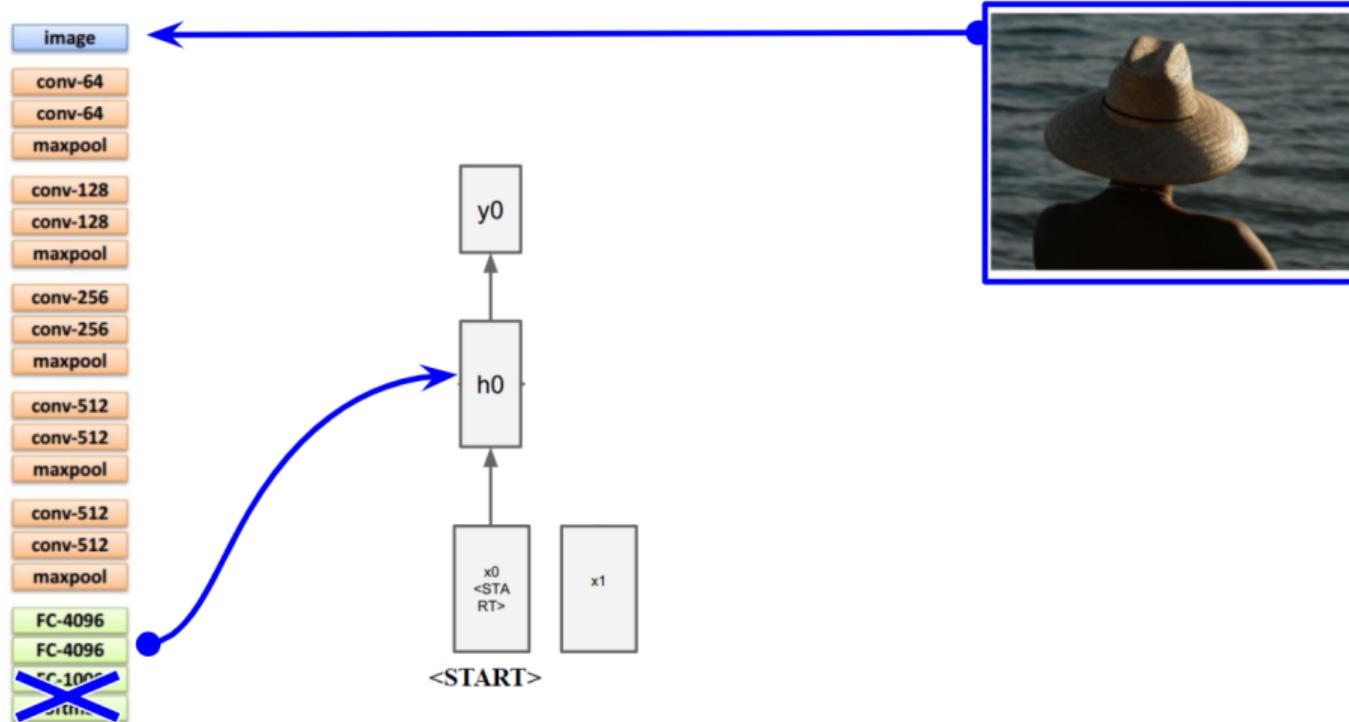
Credit: Karpathy et al, Deep visual-semantic alignments for generating image descriptions, CVPR 2015

Image Captioning: Inference (Test Time)



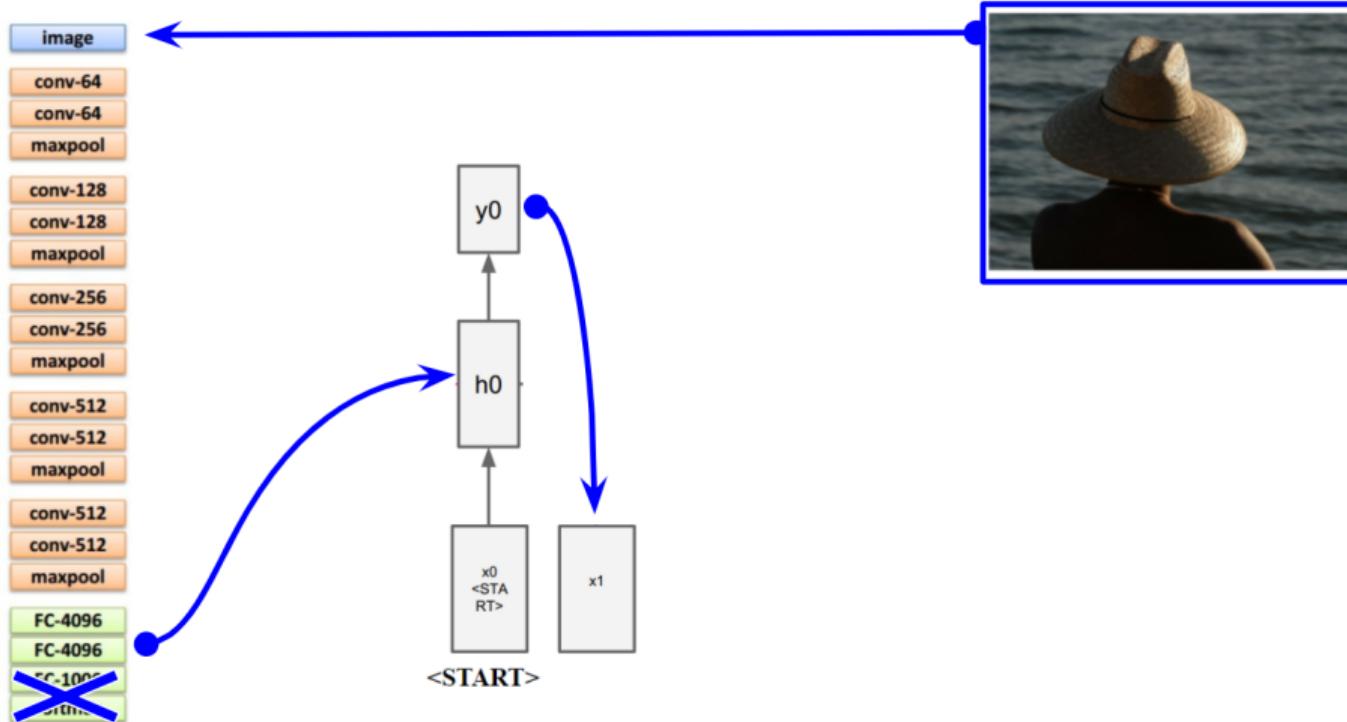
Credit: Karpathy et al, Deep visual-semantic alignments for generating image descriptions, CVPR 2015

Image Captioning: Inference (Test Time)



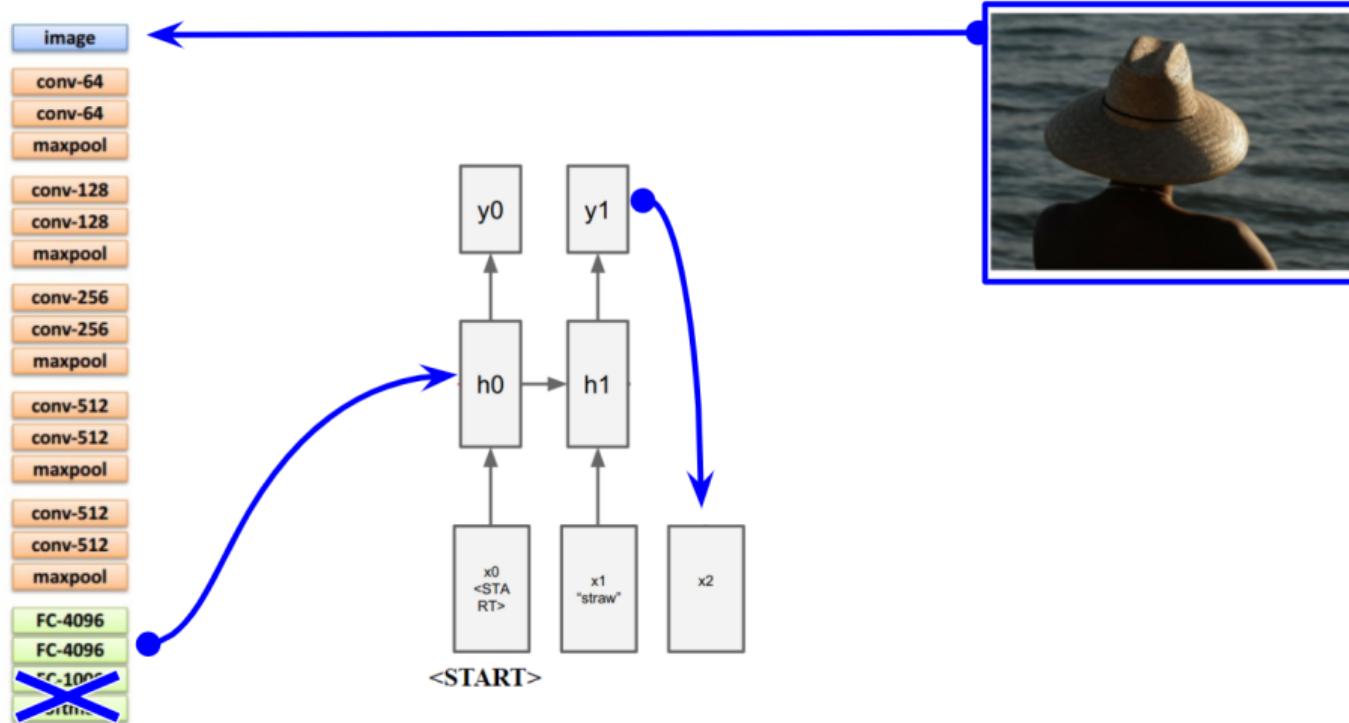
Credit: Karpathy et al, Deep visual-semantic alignments for generating image descriptions, CVPR 2015

Image Captioning: Inference (Test Time)



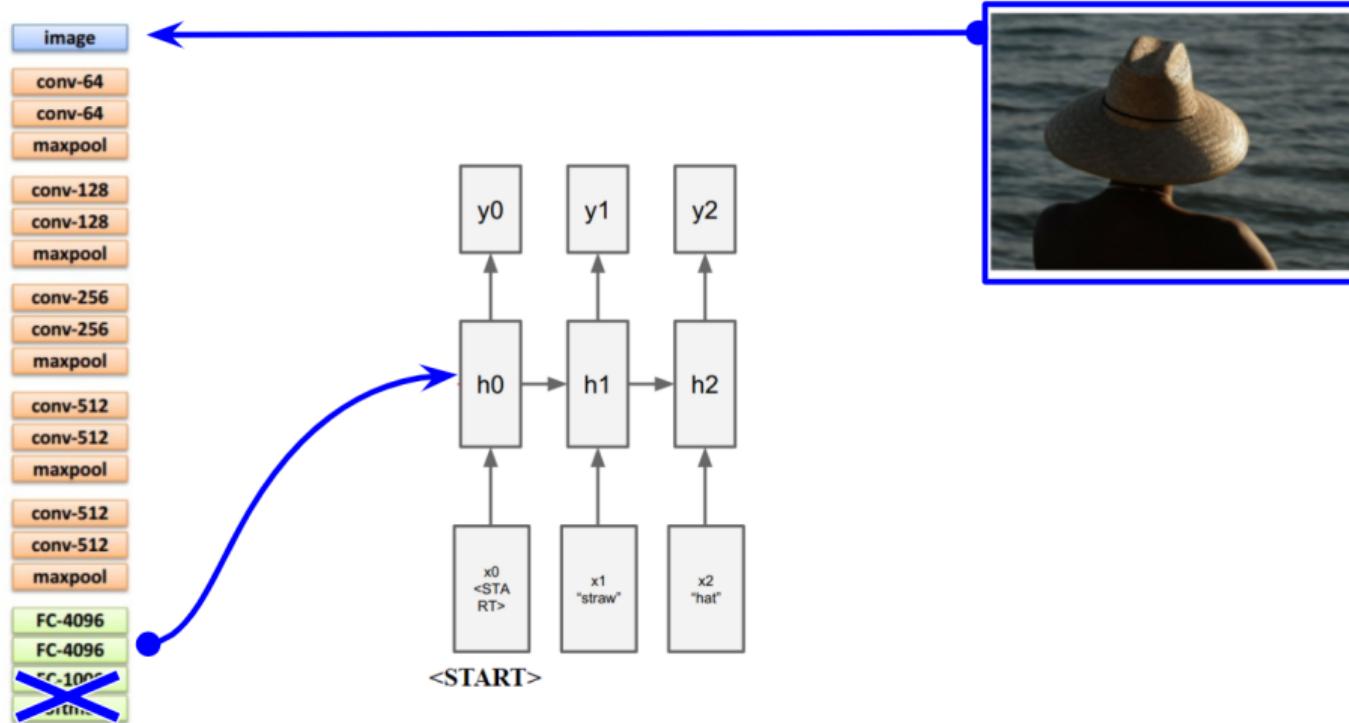
Credit: Karpathy et al, Deep visual-semantic alignments for generating image descriptions, CVPR 2015

Image Captioning: Inference (Test Time)



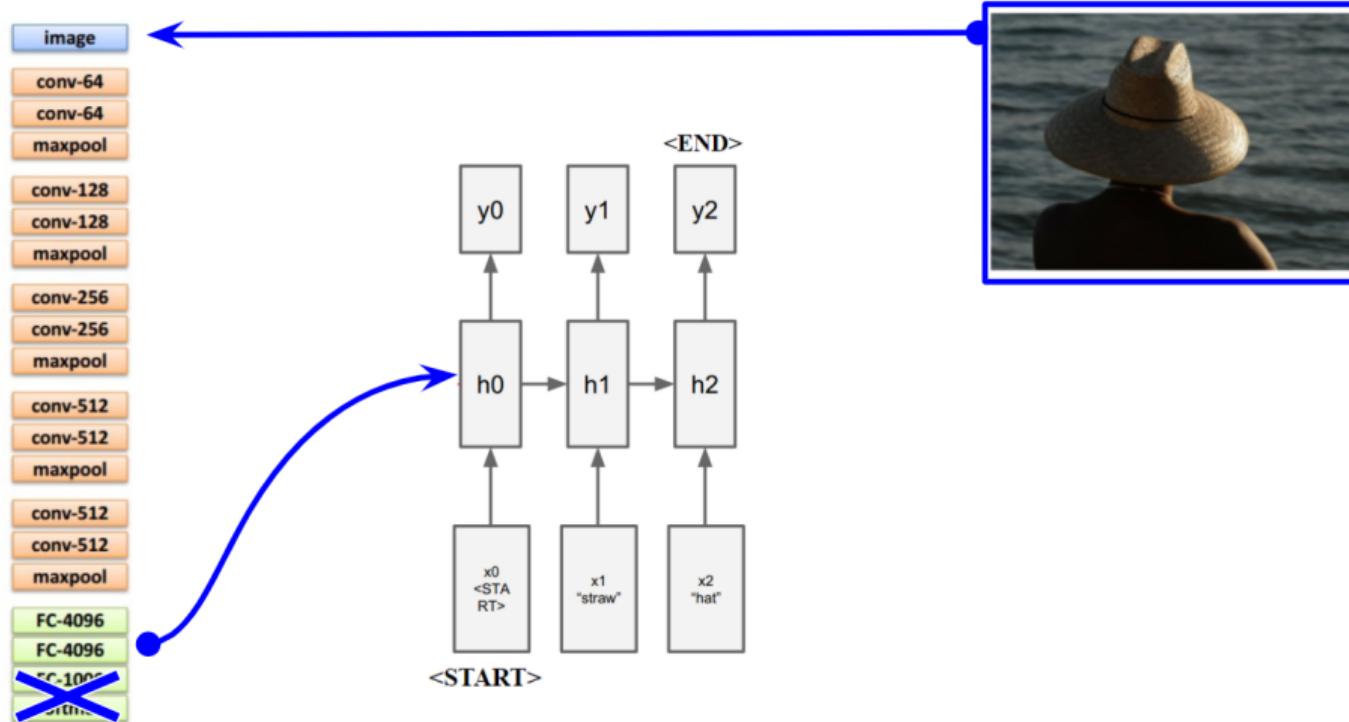
Credit: Karpathy et al, Deep visual-semantic alignments for generating image descriptions, CVPR 2015

Image Captioning: Inference (Test Time)



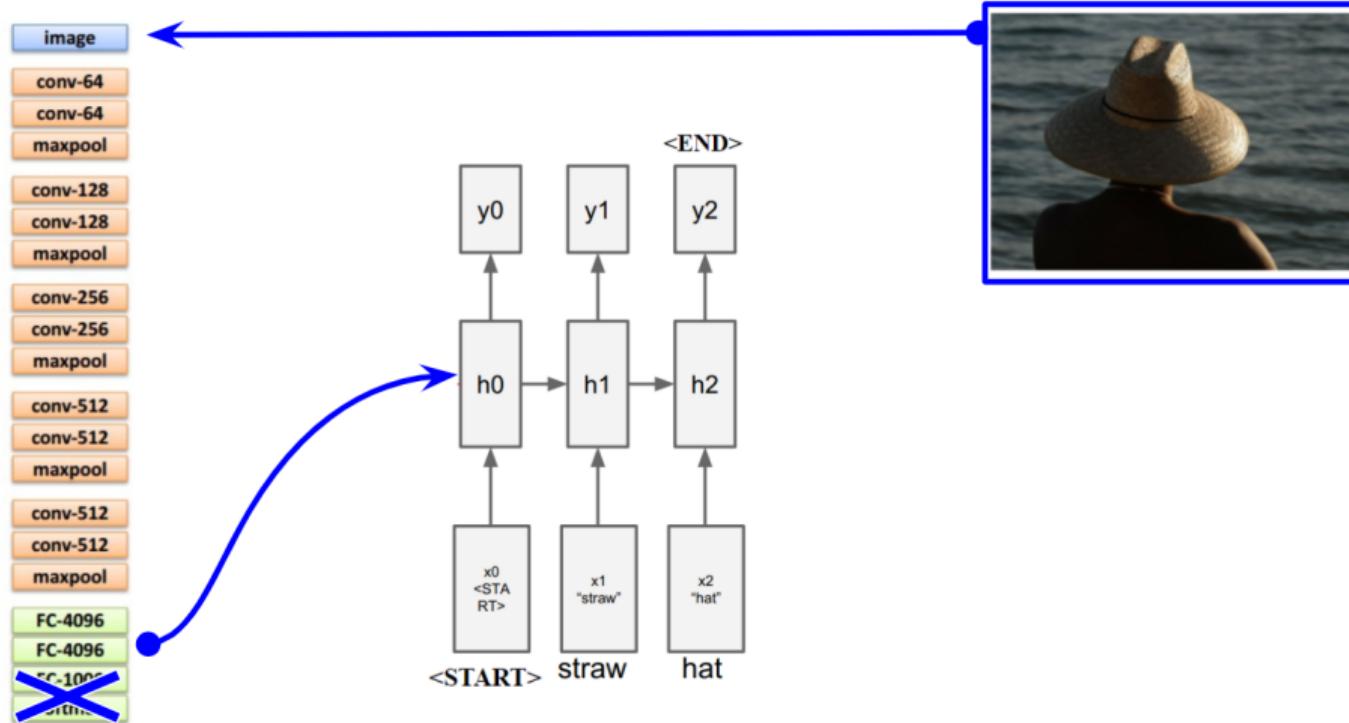
Credit: Karpathy et al, Deep visual-semantic alignments for generating image descriptions, CVPR 2015

Image Captioning: Inference (Test Time)



Credit: Karpathy et al, Deep visual-semantic alignments for generating image descriptions, CVPR 2015

Image Captioning: Inference (Test Time)



Credit: Karpathy et al, Deep visual-semantic alignments for generating image descriptions, CVPR 2015

Results



a group of people standing around a room with remotes
logprob: -9.17



a young boy is holding a baseball bat
logprob: -7.61



a cow is standing in the middle of a street
logprob: -8.84

Results: Failure Cases

Possible to understand why the method failed



a man standing next to a clock on a wall
logprob: -10.08



a young boy is holding a
baseball bat
logprob: -7.65



a cat is sitting on a couch with a remote control
logprob: -12.45

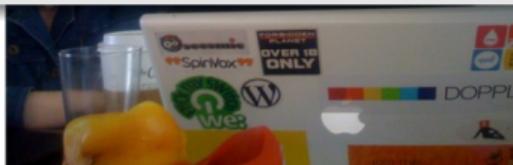
Results: Failure Cases

Not possible to understand why the method failed



How can we mitigate these failures?

Image captioning **with attention**

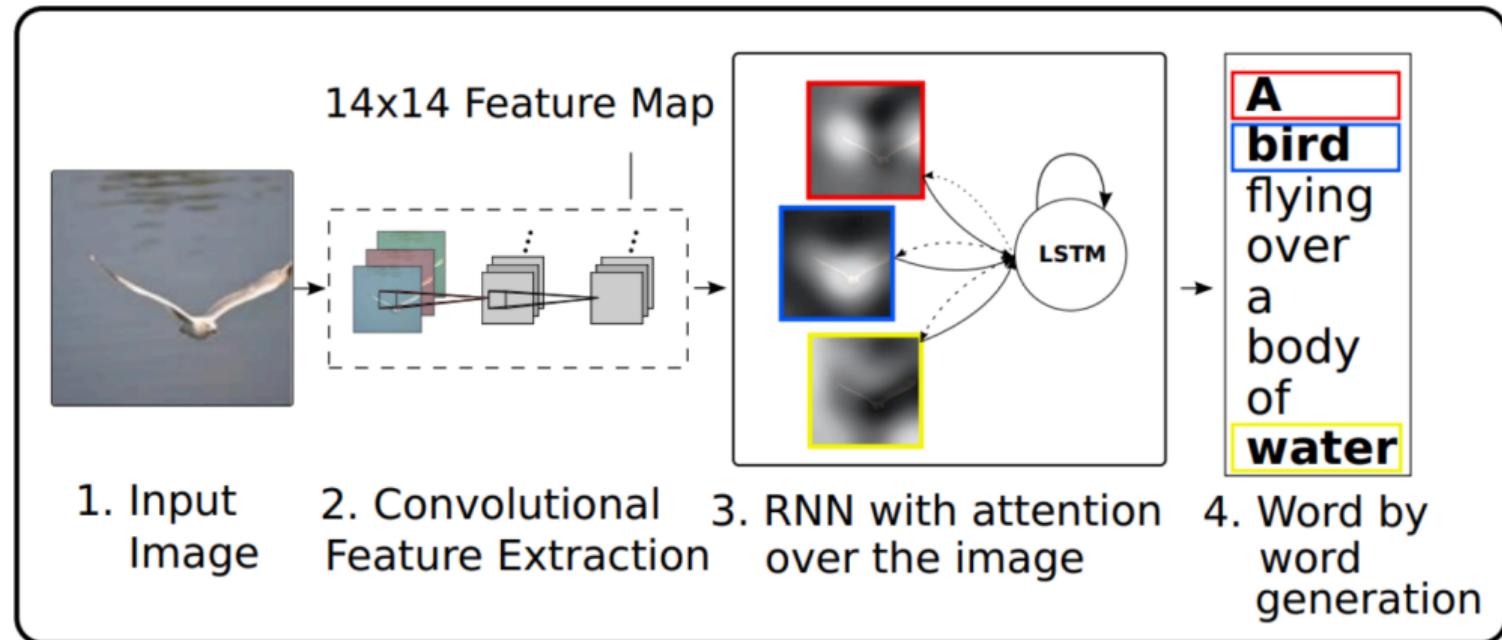


a woman holding a teddy bear in front of a mirror
logprob: -9.65



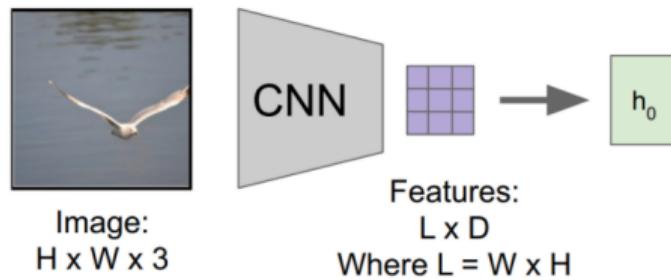
a horse is standing in the middle of a road
logprob: -10.34

Image Captioning with Attention



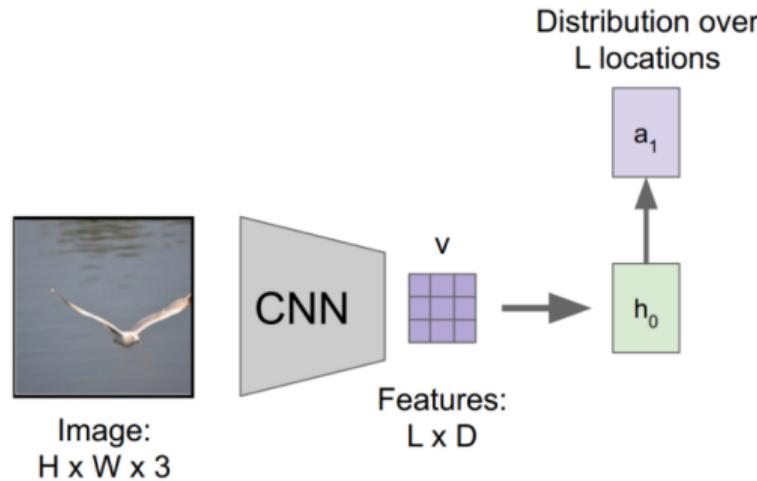
Credit: Xu et al, Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, ICML 2015

Image Captioning with Attention



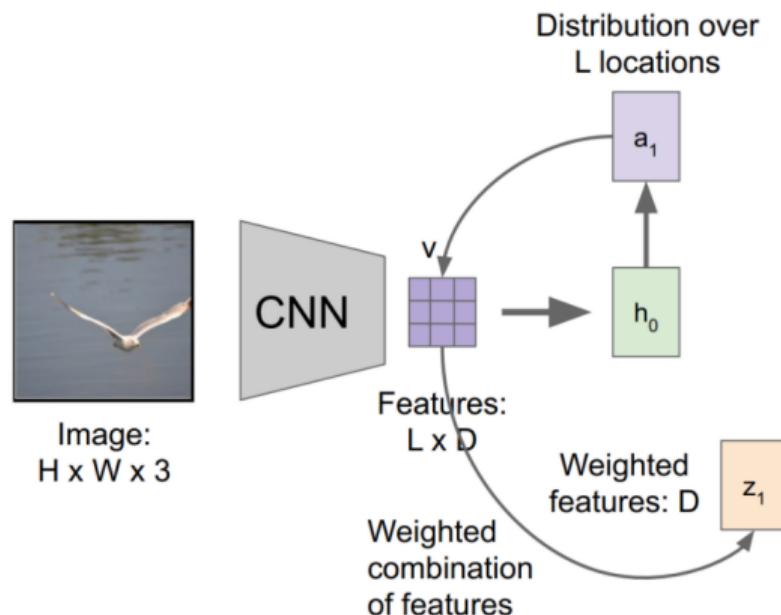
Credit: Xu et al, Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, ICML 2015

Image Captioning with Attention



Credit: Xu et al, Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, ICML 2015

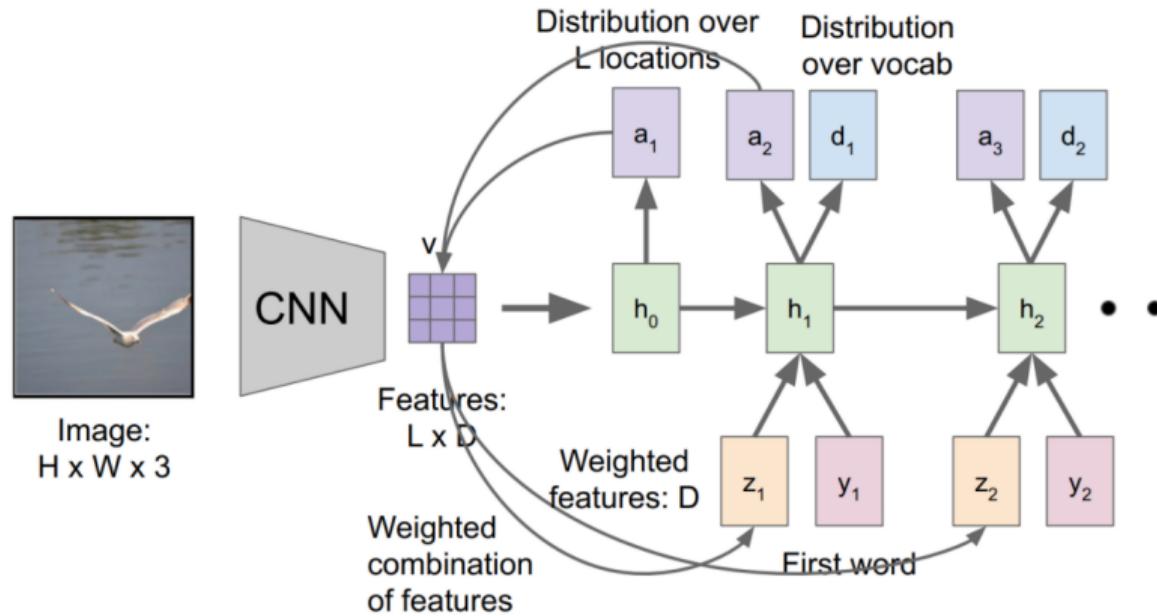
Image Captioning with Attention



$$z_1 = \sum_{i=1}^L a_i v_i$$

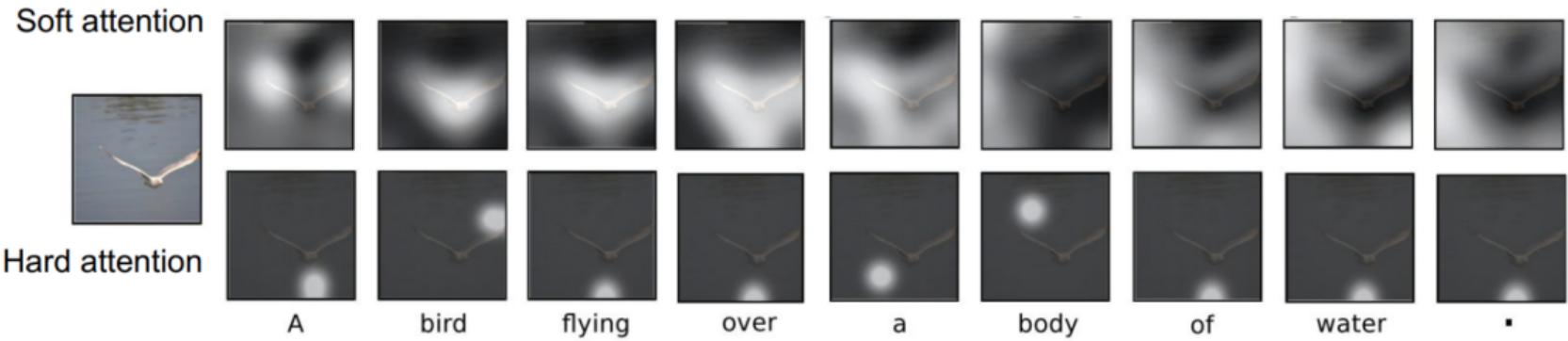
Credit: Xu et al, Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, ICML 2015

Image Captioning with Attention



Credit: Xu et al, Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, ICML 2015

Image Captioning with Attention



Credit: Xu et al, Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, ICML 2015

Image Captioning with Attention: Results



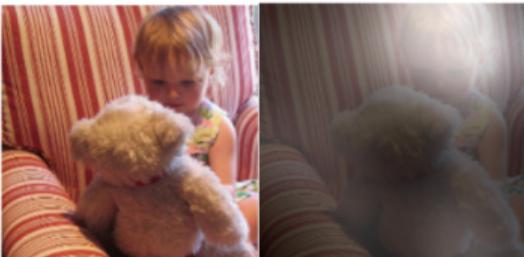
A woman is throwing a frisbee in a park.



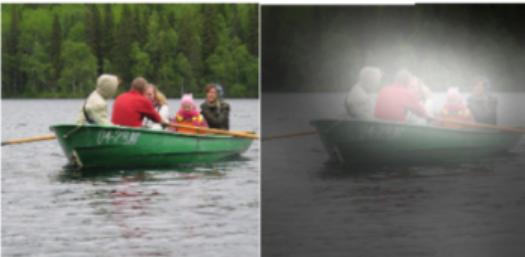
A dog is standing on a hardwood floor.



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



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



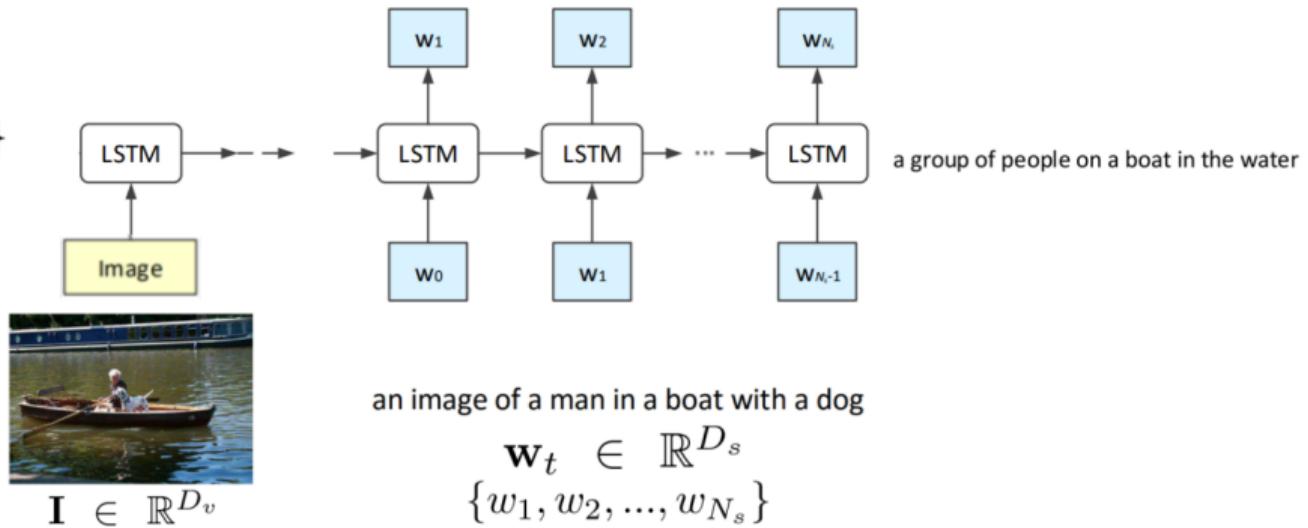
A giraffe standing in a forest with trees in the background.

Recent Efforts: Boosting Image Captioning with Attributes in LSTMs

$$\mathbf{x}^{-1} = \mathbf{T}_v \mathbf{I}, \quad \mathbf{x}^t = \mathbf{T}_s \mathbf{w}_t,$$

$$\mathbf{h}^t = f(\mathbf{x}^t),$$

$$t \in \{0, \dots, N_s - 1\}$$



Boosting Image Captioning with Attributes in LSTMs: A1

1. Transforming the images into attributes

$$\mathbf{x}^{-1} = \mathbf{T}_a \mathbf{A},$$

$$\mathbf{x}^t = \mathbf{T}_s \mathbf{w}_t,$$

$$\mathbf{h}^t = f(\mathbf{x}^t), \quad t \in \{0, \dots, N_s - 1\},$$

why? bcoz the output gonna be a caption so why not give textual data

$$\mathbf{A} \in \mathbb{R}^{D_a}$$

$$\mathcal{A} = \{a_1, a_2, \dots, a_{D_a}\}$$



Attributes:

boat: 1 water: 0.838 man: 0.762
riding: 0.728 dog: 0.547 small: 0.485
person: 0.471 river: 0.461

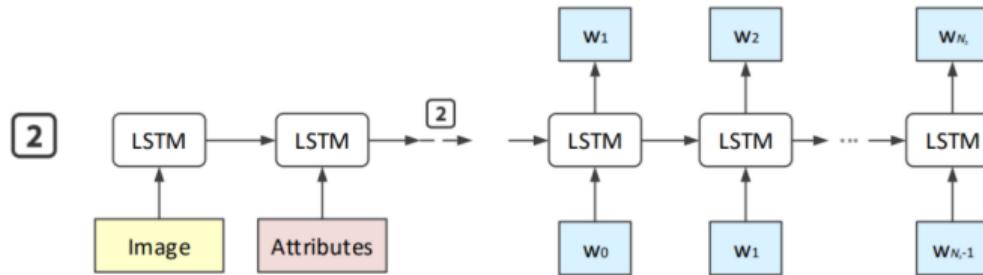
Credit: Yao et al, Boosting Image Captioning with Attributes, ICCV 2017

Boosting Image Captioning with Attributes in LSTMs: A2

$$\mathbf{x}^{-2} = \mathbf{T}_v \mathbf{I} \text{ and } \mathbf{x}^{-1} = \mathbf{T}_a \mathbf{A},$$

First: Image input
Second: Attribute input

$$\mathbf{x}^t = \mathbf{T}_s \mathbf{w}_t,$$



$$\mathbf{h}^t = f(\mathbf{x}^t), \quad t \in \{0, \dots, N_s - 1\}$$

Credit: Yao et al, Boosting Image Captioning with Attributes, ICCV 2017

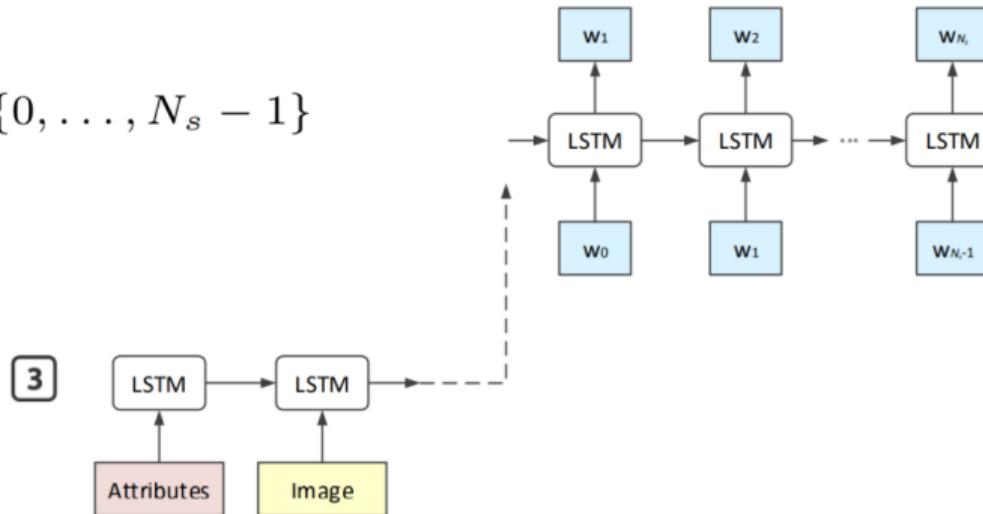
Boosting Image Captioning with Attributes in LSTMs: A3

$$\mathbf{x}^{-2} = \mathbf{T}_a \mathbf{A} \quad \text{and} \quad \mathbf{x}^{-1} = \mathbf{T}_v \mathbf{I},$$

$$\mathbf{x}^t = \mathbf{T}_s \mathbf{w}_t,$$

$$\mathbf{h}^t = f(\mathbf{x}^t), \quad t \in \{0, \dots, N_s - 1\}$$

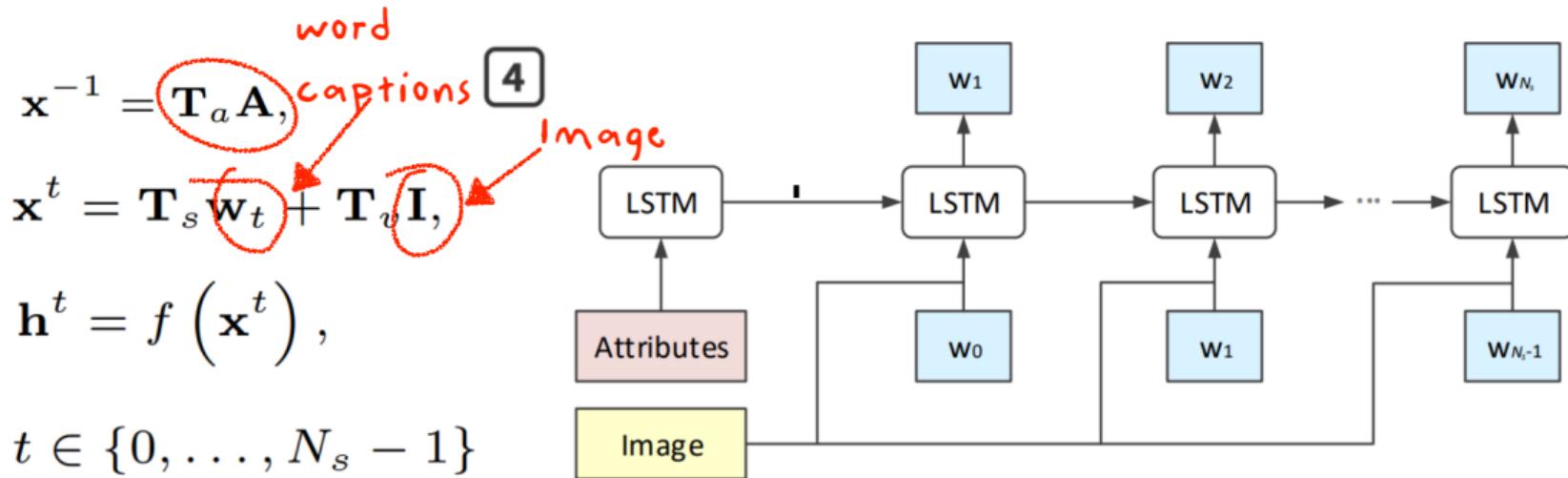
First: Attribute
Second: Image



Credit: Yao et al, Boosting Image Captioning with Attributes, ICCV 2017

Boosting Image Captioning with Attributes in LSTMs: A4

Attributes are given first only once
each time step, both image and words in caption are given as input



Credit: Yao et al, Boosting Image Captioning with Attributes, ICCV 2017

Boosting Image Captioning with Attributes in LSTMs: A5

opposite of A4

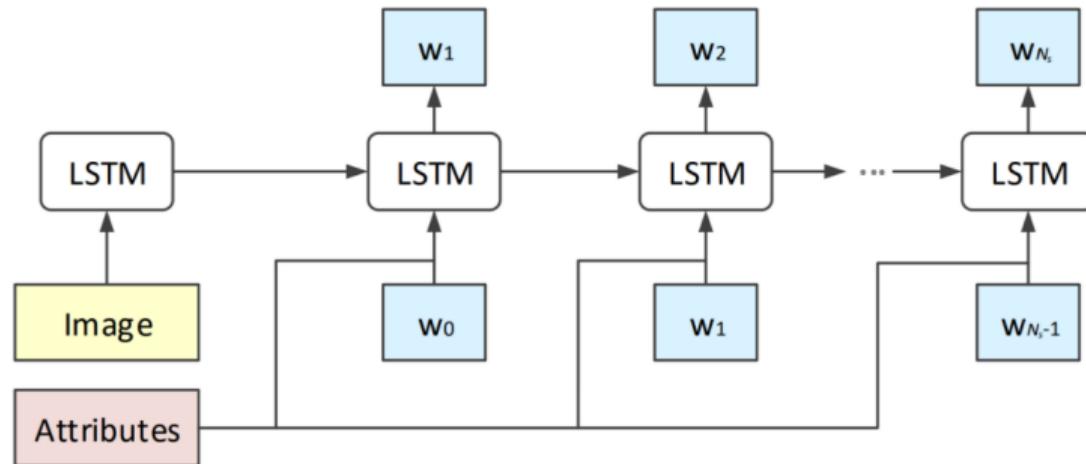
$$\mathbf{x}^{-1} = \mathbf{T}_v \mathbf{I},$$

5

$$\mathbf{x}^t = \mathbf{T}_s \mathbf{w}_t + \mathbf{T}_a \mathbf{A},$$

$$\mathbf{h}^t = f(\mathbf{x}^t),$$

$$t \in \{0, \dots, N_s - 1\}$$

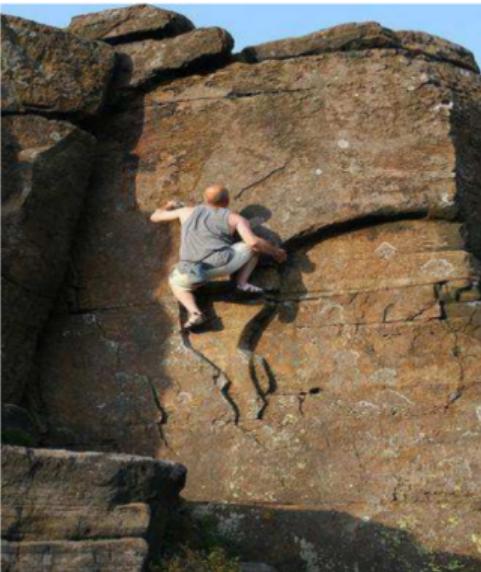


Credit: Yao et al, Boosting Image Captioning with Attributes, ICCV 2017

Boosting Image Captioning with Attributes in LSTMs: Observations

- $\text{LSTM-}A_1 > \text{LSTM}$
 - Indicates advantage of exploiting high-level attributes than image representations
- $\text{LSTM-}A_2 > \text{LSTM-}A_1$
 - Integrating image representations performs better
- $\text{LSTM-}A_3 > \text{LSTM-}A_2$
 - Benefits from mechanism of first feeding high-level attributes into LSTM instead of starting from image representations
- $\text{LSTM-}A_4 < \text{LSTM-}A_3$
 - This may be because noise in image can be explicitly accumulated, and thus network overfits more easily
 - But $\text{LSTM-}A_5$ which feeds attributes at each time step shows improvements on $\text{LSTM-}A_3$
- $\text{LSTM-}A_2, \text{LSTM-}A_3, \text{LSTM-}A_5 > \text{LSTM}$
 - Indicates that image representations and attributes are complementary and have mutual reinforcement for image captioning

StyleNet: Generating Attractive Visual Captions with Styles



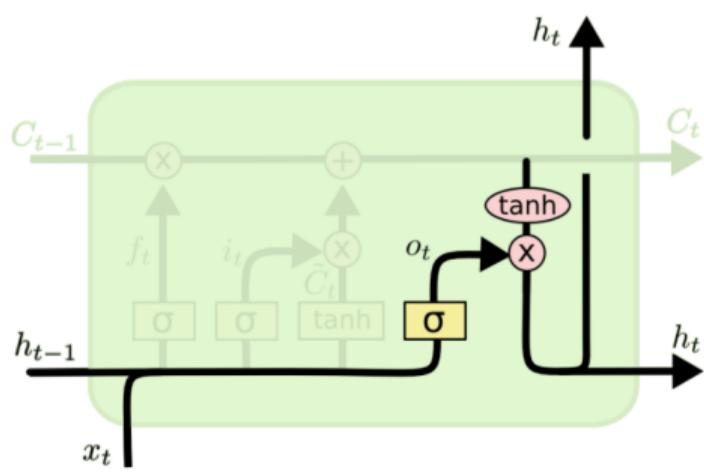
CaptionBot: A man on a rocky hillside next to a stone wall.

Romantic: A man uses rock climbing to conquer the high.

Humorous: A man is climbing the rock like a lizard.

StyleNet: Generating Attractive Visual Captions with Styles

StyleNet proposes a factored LSTM



$$W_x = U_x S_x V_x$$

$$i_t = \text{sigmoid}(W_{ix}x_t + W_{ih}h_{t-1})$$

$$f_t = \text{sigmoid}(W_{fx}x_t + W_{fh}h_{t-1})$$

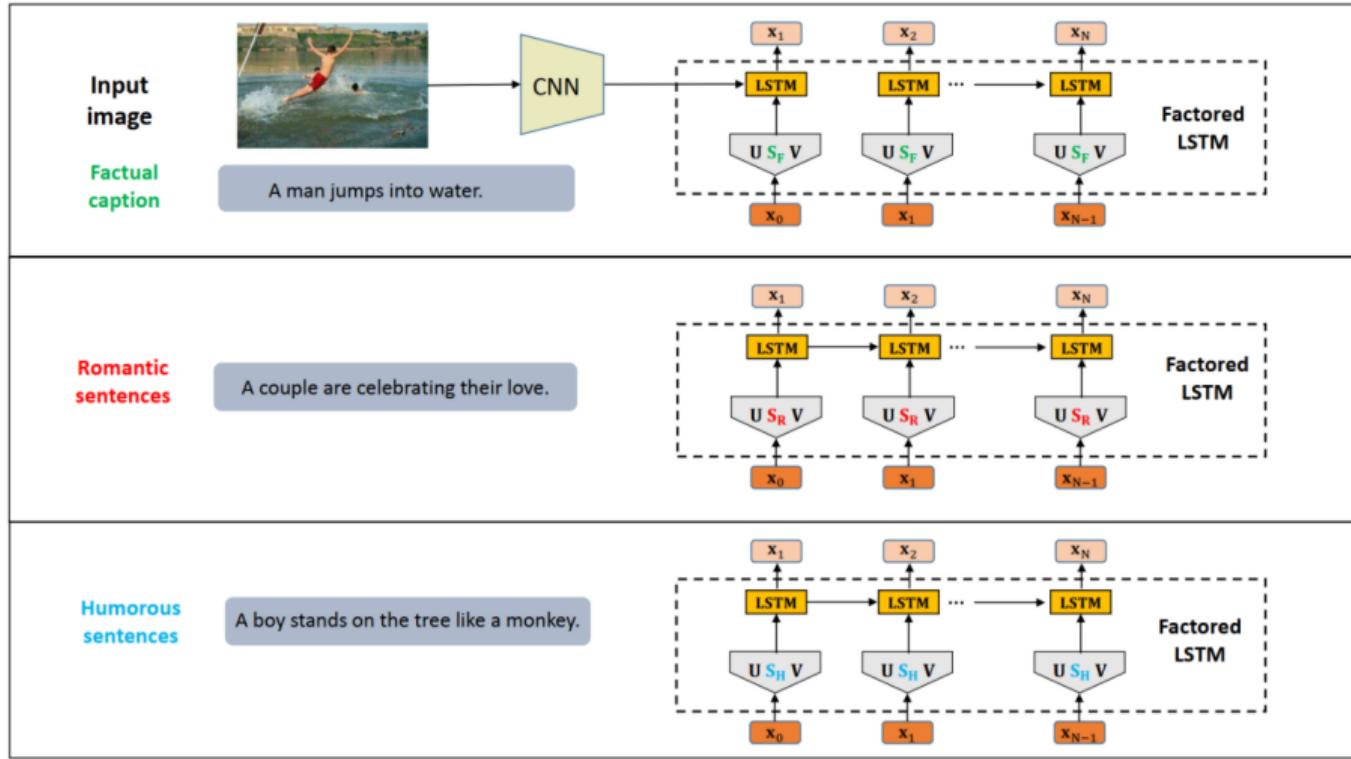
$$o_t = \text{sigmoid}(W_{ox}x_t + W_{oh}h_{t-1})$$

$$\tilde{c}_t = \tanh(W_{cx}x_t + W_{ch}h_{t-1})$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$

$$h_t = o_t \odot h_t$$

StyleNet: Generating Attractive Visual Captions with Styles



StyleNet: Generating Attractive Visual Captions with Styles



F: A snowboarder in the air .

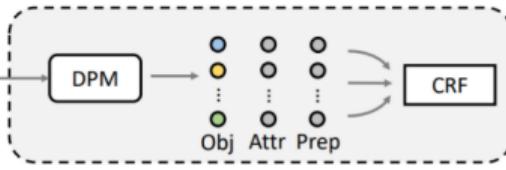
R: A man is doing a trick on a skateboard to show his courage .

H: A man is jumping on a snowboard to reach outer space .

At test swap S_x accordingly to get the desired output type.

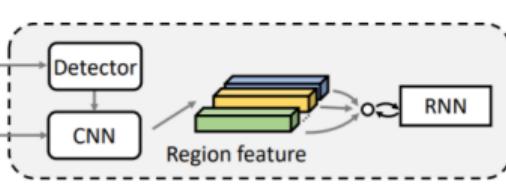
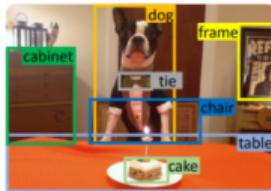
Neural Baby Talk

More Grounded



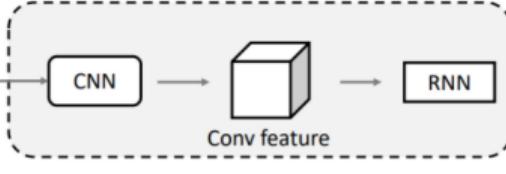
This is a photograph
of one dog and one
cake. The dog is ...

Neural Baby Talk



A (yellow) with a (grey) is
sitting at (blue) with
a (green).
— puppy — tie
— cake — table

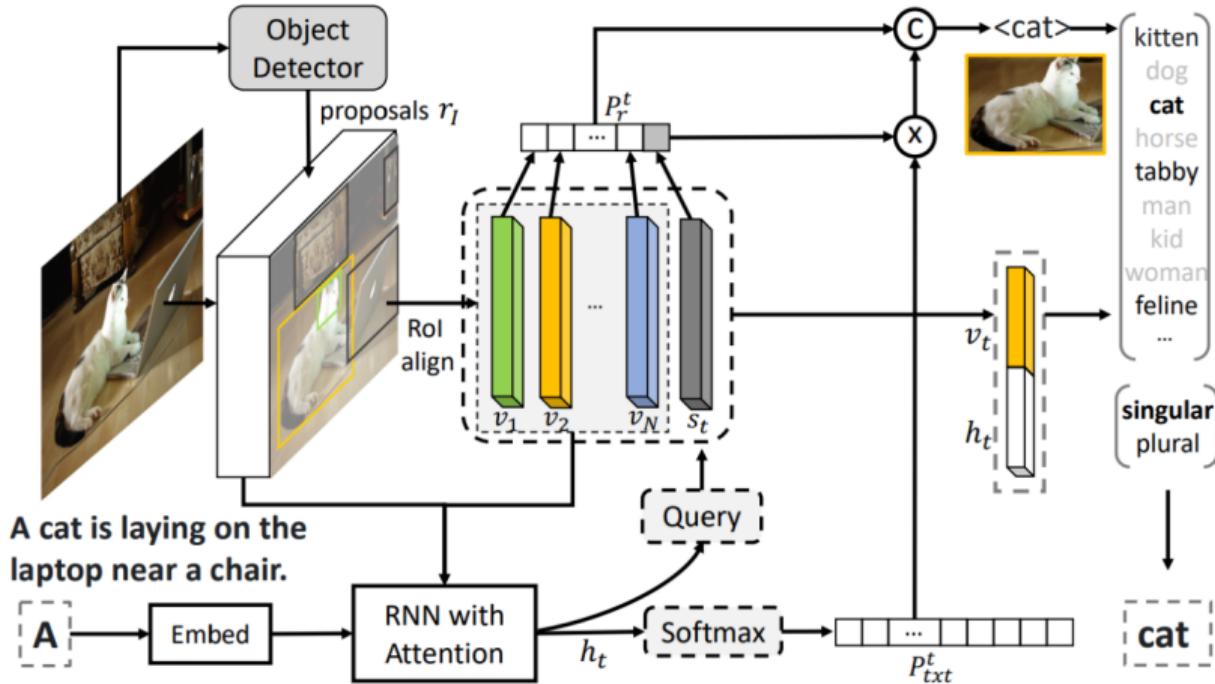
More Natural



A dog is sitting on a
couch with a toy.

Credit: Lu et al, Neural Baby Talk, CVPR 2018

Neural Baby Talk



Credit: Lu et al, Neural Baby Talk, CVPR 2018