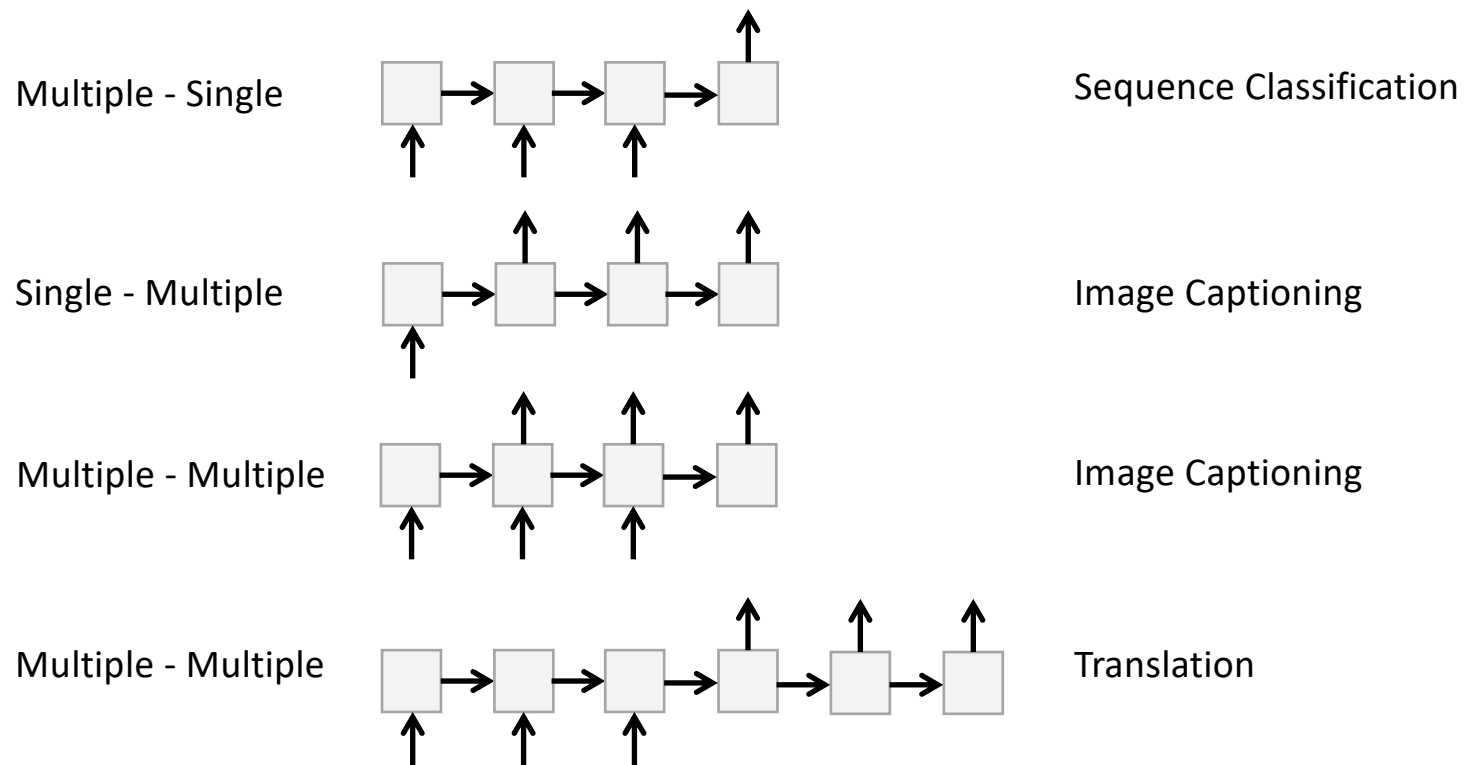


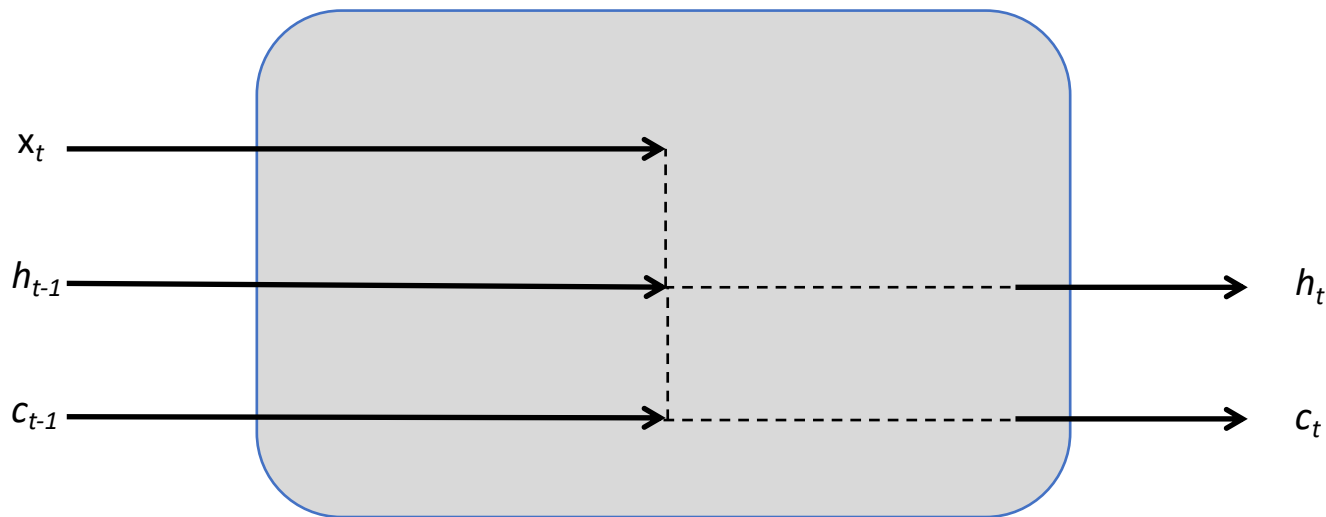
# Sequence models





# Long Short-Term Memory (LSTM)

- Add a *memory cell* that is not subject to matrix multiplication or squishing, thereby avoiding gradient decay



- RNN harder to train – more specifics
- For many of the tasks temporal models are
- Pros of LSTMs
  - - variable input sequences
  - - variable output sequences
- Hard to train – problems with modelling longer sequences (beyond ~ 100 steps)
  - both the gate, memory cell models

*An Empirical Evaluation of Generic Convolutional and Recurrent Networks for Sequence Modeling, Bai et al., arXiv 2018*



# Temporal models from processing video

- Relevant for:
- Action recognition
- Tracking
- Monocular motion estimation
- Future prediction

1. Keep frames independent
2. Temporal model fit RNN

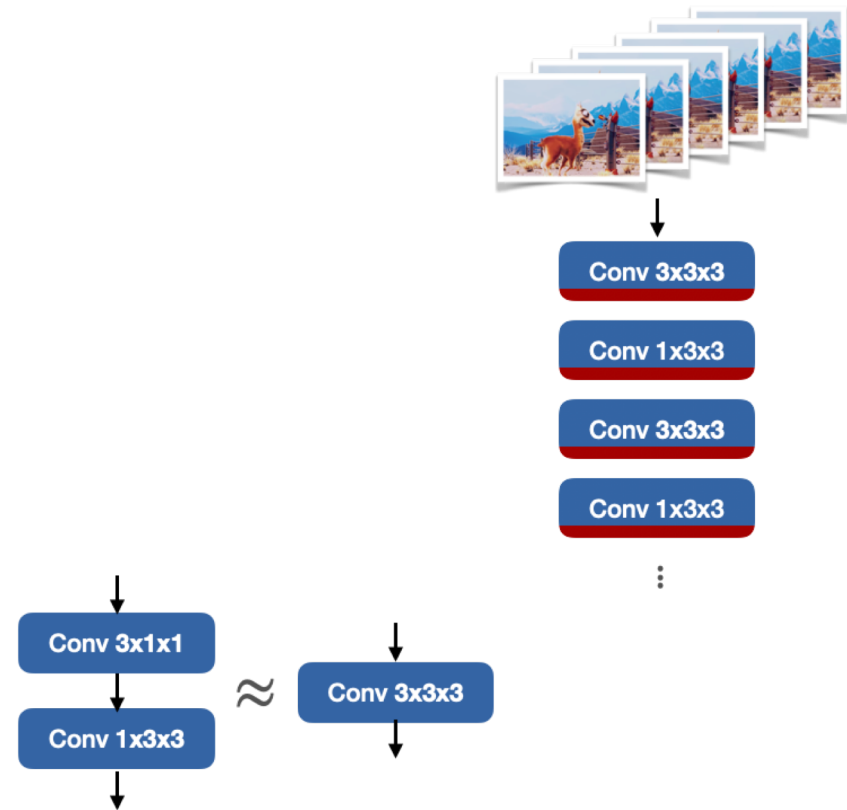


# Temporal convolutions

- 3D CNN
- Add temporal dimension to convolution
- 3D kernel
- Slide the kernel through space and time
- $w \times h \times 3 \times t$
- Too many parameters, too slow
- Idea separate spatial and temporal convolutions ( $C$  is number of channels)

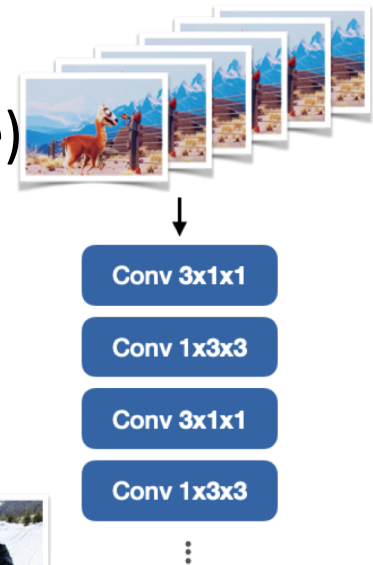
$$O((WH+T) \times C^2)$$

$$O(WHT) \times C^2$$



# Temporal models – Activity Recognition

- Train representation on Image Net 2D CNN
- Inflate the network in time (replicated the kernel in time)
- Or insert 1D temporal convolutions initialized by (averaging and fine tune on video dataset), faster training still computationally demanding



# Computer Vision Problems

- Activities of daily living
- Relationship between space and time
- Representations obtained by vision algorithms often do not, results imperfect and often not all relevant For the task at hand, e.g. Autonomous Driving

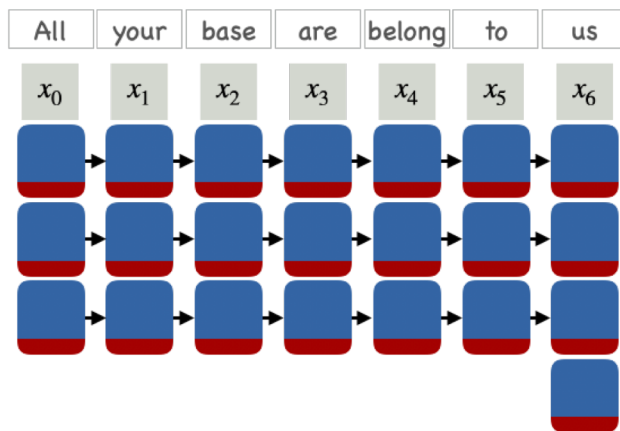


Image source: EPIC-Kitchens dataset, Damen et al., <https://arxiv.org/abs/1804.02748>

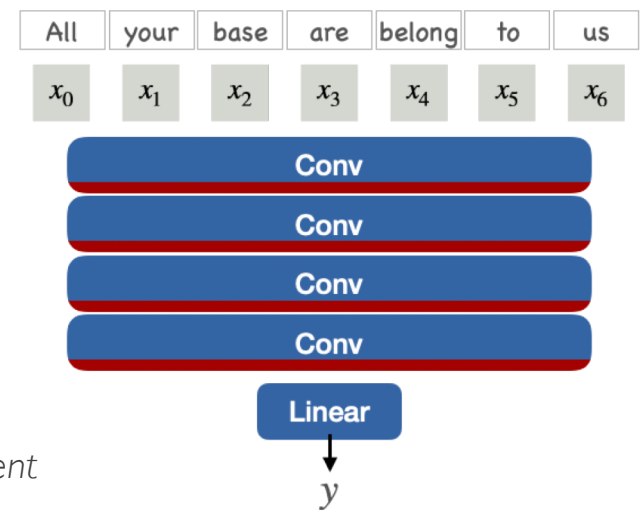


# Temporal models in general

- Recurrent models – input sequence output – predicts output – for next word prediction task
- We can handle the variable length – use convolutions (dilated convolutions – the receptive field grows with the number of layers)
- How to handle temporal models for the same task



*An Empirical Evaluation of Generic Convolutional and Recurrent Networks for Sequence Modeling, Bai et al., arXiv 2018*

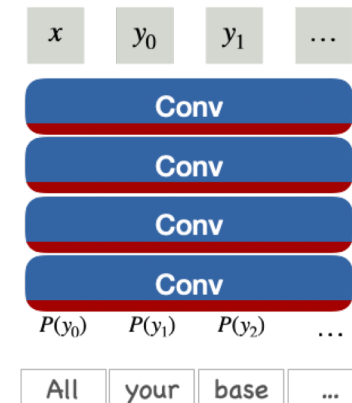
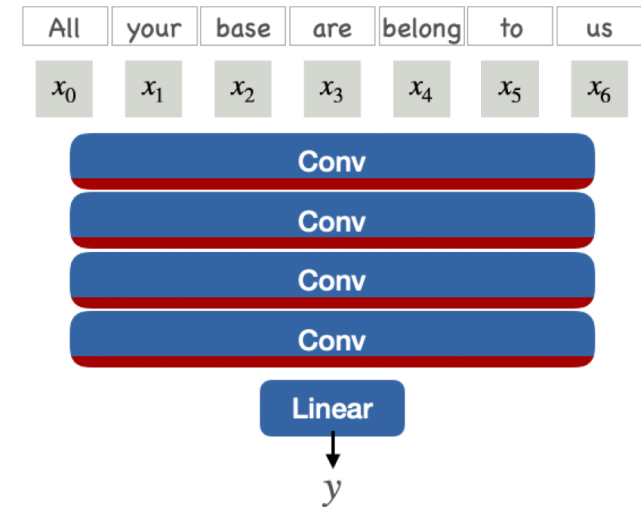
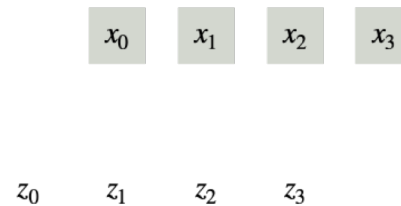


# Temporal convolutions

- Dilated convolutions
- 5-10 layers – 100 steps context
- Caveat - need causal convolutions
- Only look in the past (works in 1D)
- Autoregressive model

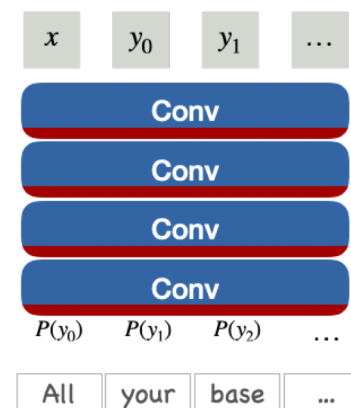
$$P(y_0 | x) \cdot P(y_1 | x, y_0) \cdot P(y_2 | x, y_0, y_1) \cdot \dots$$

- Shift input



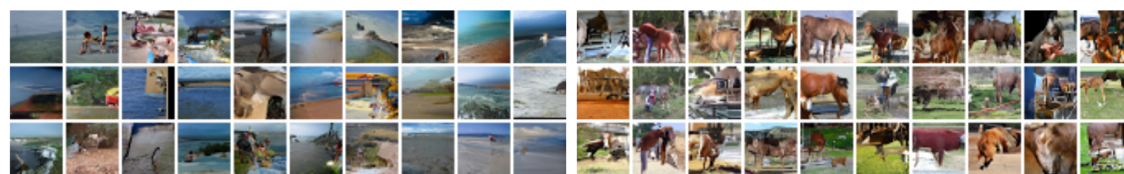
- How to generate sequences ?
- Need to modify convolutions to look in the past not the future, harder for images

New type of decoder conditioned on image encoder



African elephant

Coral Reef

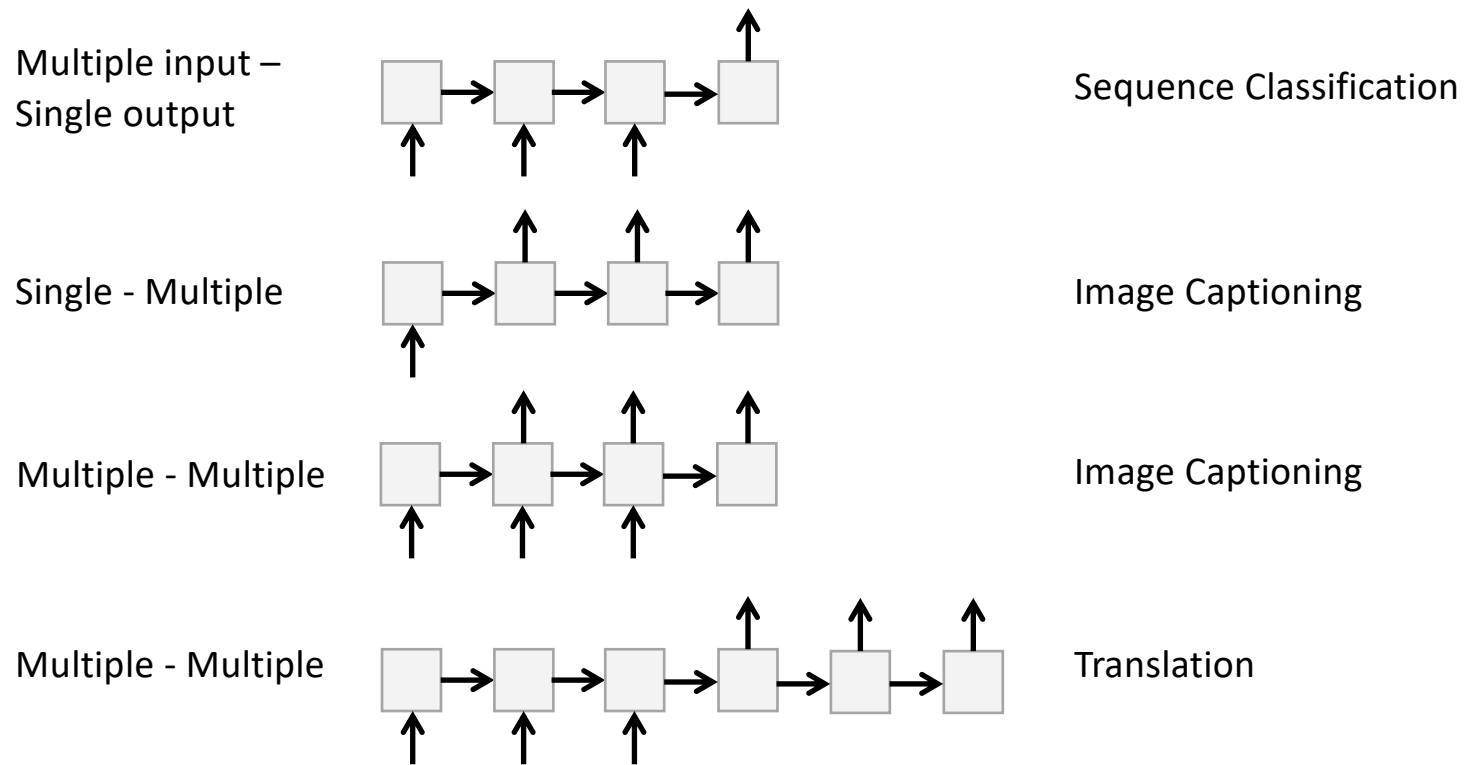


Sandbar

Sorrel horse

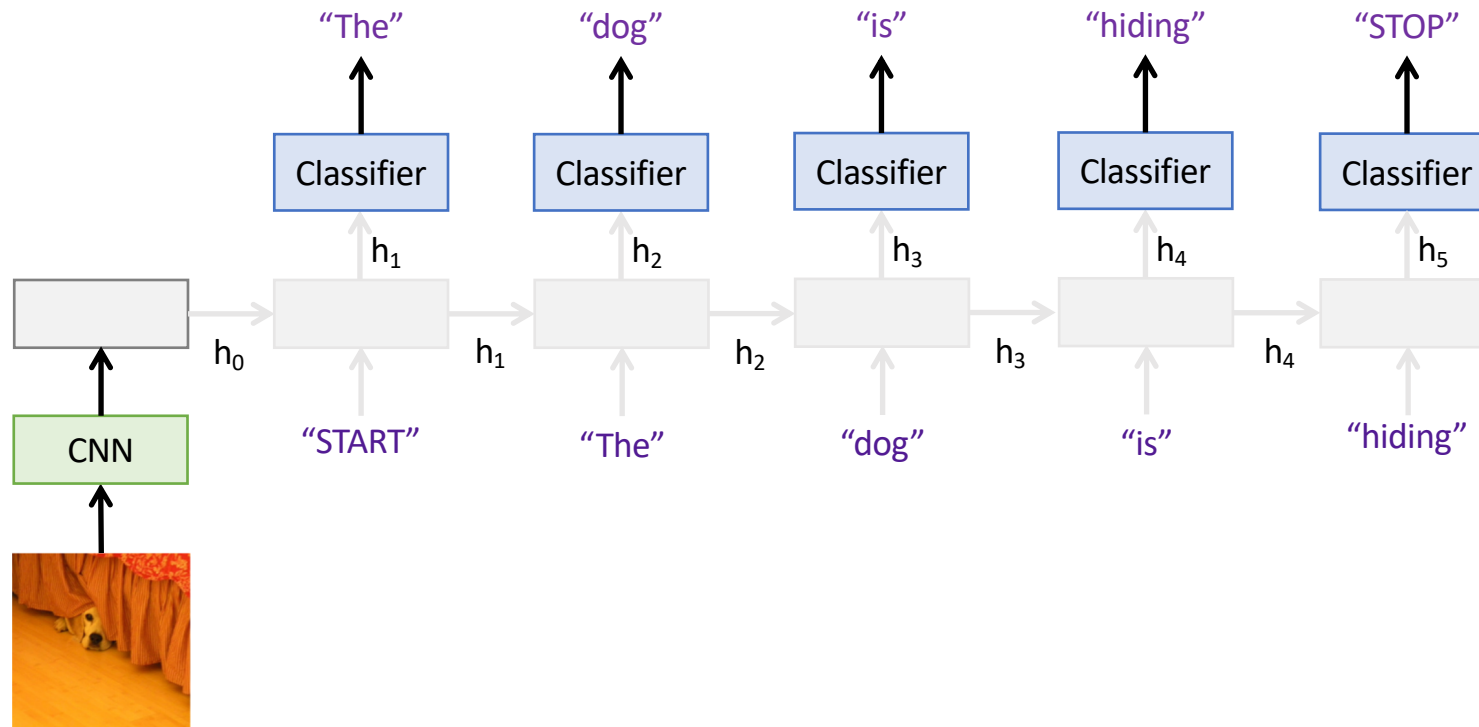
Conditional image generation with pixelCNN decoders, [Aaron van den Oord](#), [Nal Kalchbrenner](#), [Oriol Vinyals](#), [Lasse Espeholt](#), [Alex Graves](#), [Koray Kavukcuoglu](#)

# RNN use Cases

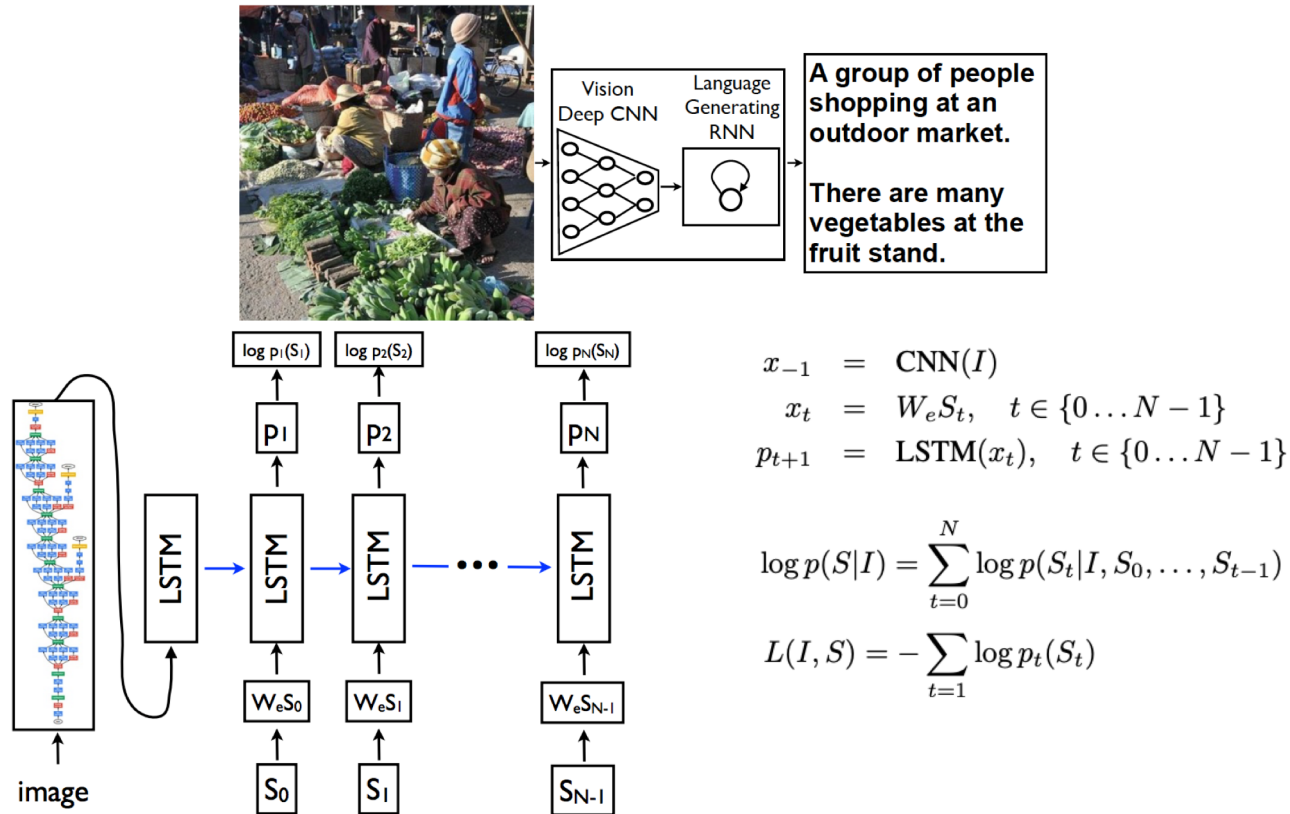




# Review: Image captioning

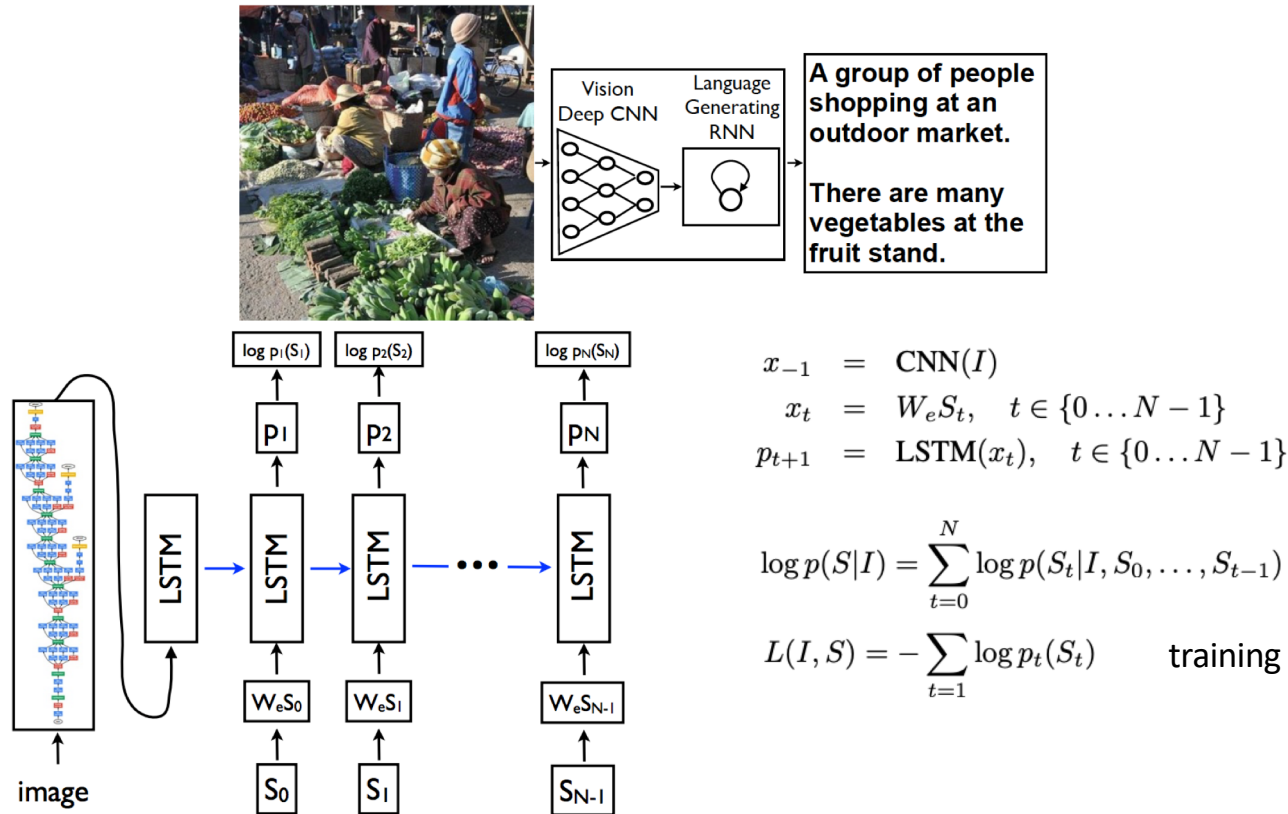


# Image Caption Generation



O. Vinyals, A. Toshev, S. Bengio, D. Erhan, [Show and Tell: A Neural Image Caption Generator](#), CVPR 2015

# Image Caption Generation



O. Vinyals, A. Toshev, S. Bengio, D. Erhan, [Show and Tell: A Neural Image Caption Generator](#), CVPR 2015

# Image Caption Generation

A person riding a motorcycle on a dirt road.



Two dogs play in the grass.



A skateboarder does a trick on a ramp.



A dog is jumping to catch a frisbee.



A group of young people playing a game of frisbee.



Two hockey players are fighting over the puck.



A little girl in a pink hat is blowing bubbles.



A refrigerator filled with lots of food and drinks.



A herd of elephants walking across a dry grass field.



A close up of a cat laying on a couch.



A red motorcycle parked on the side of the road.



A yellow school bus parked in a parking lot.



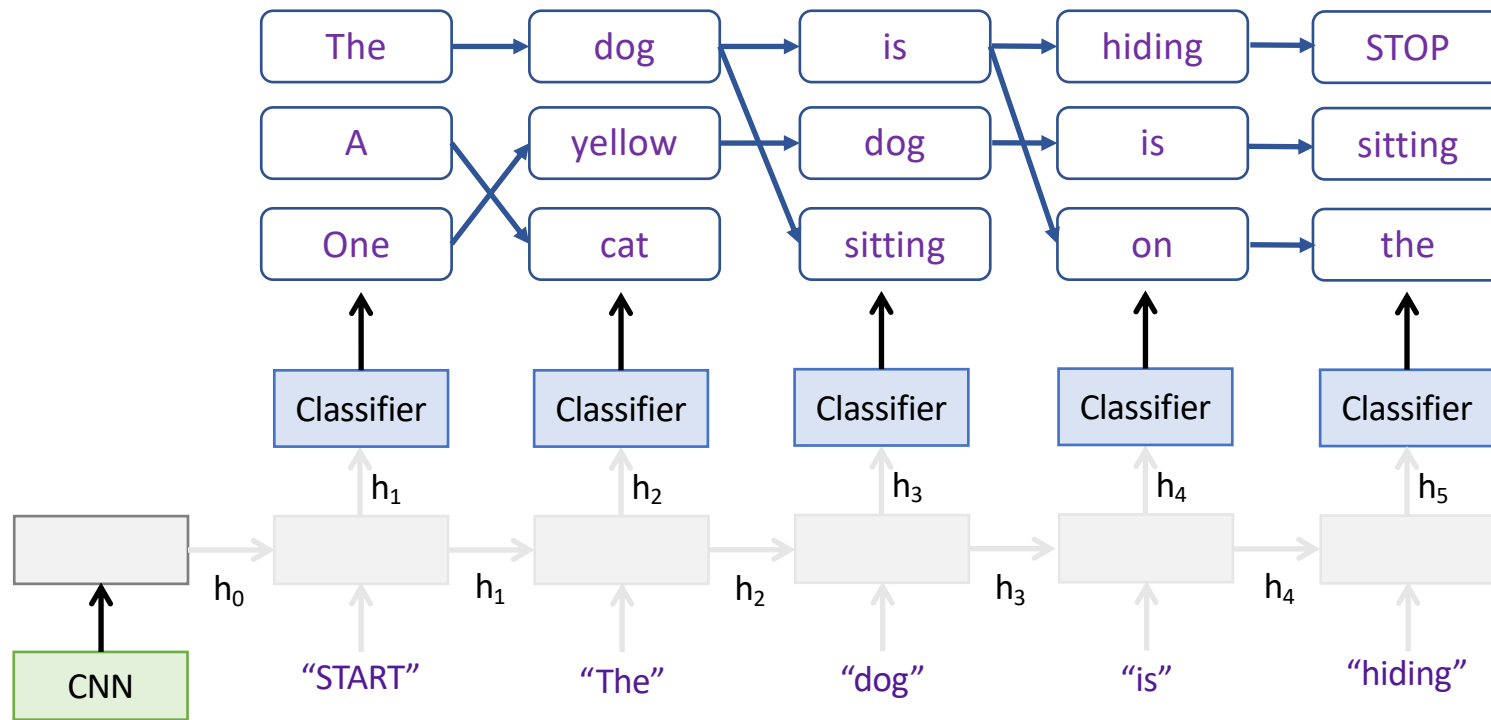
Describes without errors

Describes with minor errors

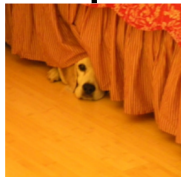
Somewhat related to the image

Unrelated to the image

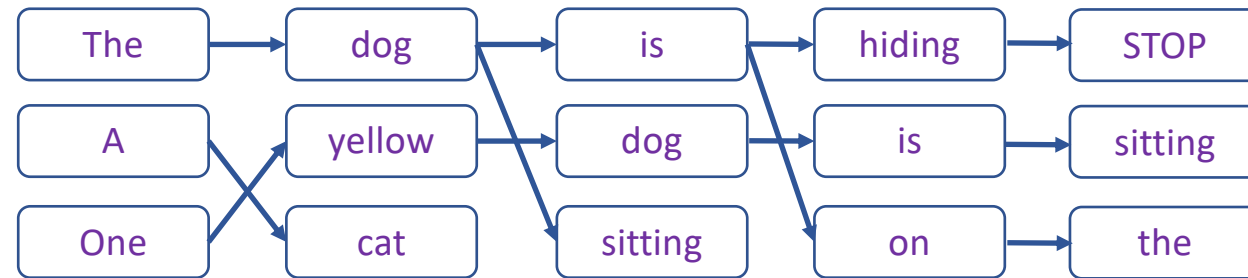
# Review: Image captioning



Different strategies how to generate the final caption:  
Sampling : sample the first word based on probability  $p_1$  and use its embedding at input  
sample the second word etc.



# Beam search



- Maintain  $k$  top-scoring candidate sentences (according to sum of per-word log-likelihoods)
  - At each step, generate all their successors and reduce to  $k$  (*beam width*)

# How to evaluate image captioning?



- Reference sentences (written by human annotators):
  - “A dog hides underneath a bed with its face peeking out of the bed skirt”
  - “The small white dog is peeking out from under the bed”
  - “A dog is peeking its head out from underneath a bed skirt”
  - “A dog peeking out from under a bed”
  - “A dog that is under a bed on the floor”
- Generated sentence:
  - “A dog is hiding”



# BLEU: Bilingual Evaluation Understudy

- **N-gram precision:** count the number of n-gram matches between candidate and reference translation, divide by total number of n-grams in candidate translation
  - Clip counts by the maximum number of times an n-gram occurs in any reference translation
  - Multiply by *brevity penalty* to penalize short translations
- Most commonly used measure despite well-known shortcomings
- $H(i)$  number of i-gram tuples,  $Matched(i)$  is number of times tuple occurs in hypothesis, min with number of times tuple occurs in reference

$$P(i) = \frac{Matched(i)}{H(i)} \qquad Matched(i) = \sum_{t_i} \min\{C_h(t_i), \max_j C_{h_j}(t_i)\}$$

$$BLEU_a = \left\{ \prod_{i=1}^N P(i) \right\}^{1/N}$$

K. Papineni, S. Roukos, T. Ward, W.-J. Zhu, [BLEU: a Method for Automatic Evaluation of Machine Translation](#), ACL 2002





- Overview
- Challenges
- Download
- Evaluate
- Leaderboard

Table-C5

Table-C40

2015 Captioning Challenge

Last update: June 8, 2015. Visit [CodaLab](#) for the latest results.

	CIDEr-D	Meteor	ROUGE-L	BLEU-1	BLEU-2	BLEU-3	BLEU-4	
m-RNN (Baidu/ UCLA) <sup>[16]</sup>	0.886	0.238	0.524	0.72	0.553	0.41	0.302	
m-RNN <sup>[15]</sup>	0.817	0.218	0.501	0.718	0.515	0.401	0.299	
MSR Captiva							0.308	
Google <sup>[4]</sup>	CIDEr-D	CIDEr: Consensus-based Image Description Evaluation						0.309
Berkeley LR	METEOR	Meteor Universal: Language Specific Translation Evaluation for Any Target Language						0.277
Nearest Neig	Rouge-L	ROUGE: A Package for Automatic Evaluation of Summaries						0.28
MSR <sup>[8]</sup>	BLEU	BLEU: a Method for Automatic Evaluation of Machine Translation						0.291
Montreal/Toronto <sup>[10]</sup>	0.85	0.243	0.513	0.689	0.515	0.372	0.268	
PicSOM <sup>[13]</sup>	0.833	0.231	0.505	0.683	0.51	0.377	0.281	
Tsinghua Bigeye <sup>[14]</sup>	0.673	0.207	0.49	0.671	0.494	0.35	0.241	
MLBL <sup>[7]</sup>	0.74	0.219	0.499	0.666	0.498	0.362	0.26	
Human <sup>[5]</sup>	0.854	0.252	0.484	0.663	0.469	0.321	0.217	

Metrics

CIDEr-D

CIDEr: Consensus-based Image Description Evaluation

METEOR

Meteor Universal: Language Specific Translation Evaluation for Any Target Language

Rouge-L

ROUGE: A Package for Automatic Evaluation of Summaries

BLEU

BLEU: a Method for Automatic Evaluation of Machine Translation

<http://mscoco.org/dataset/#captions-leaderboard>



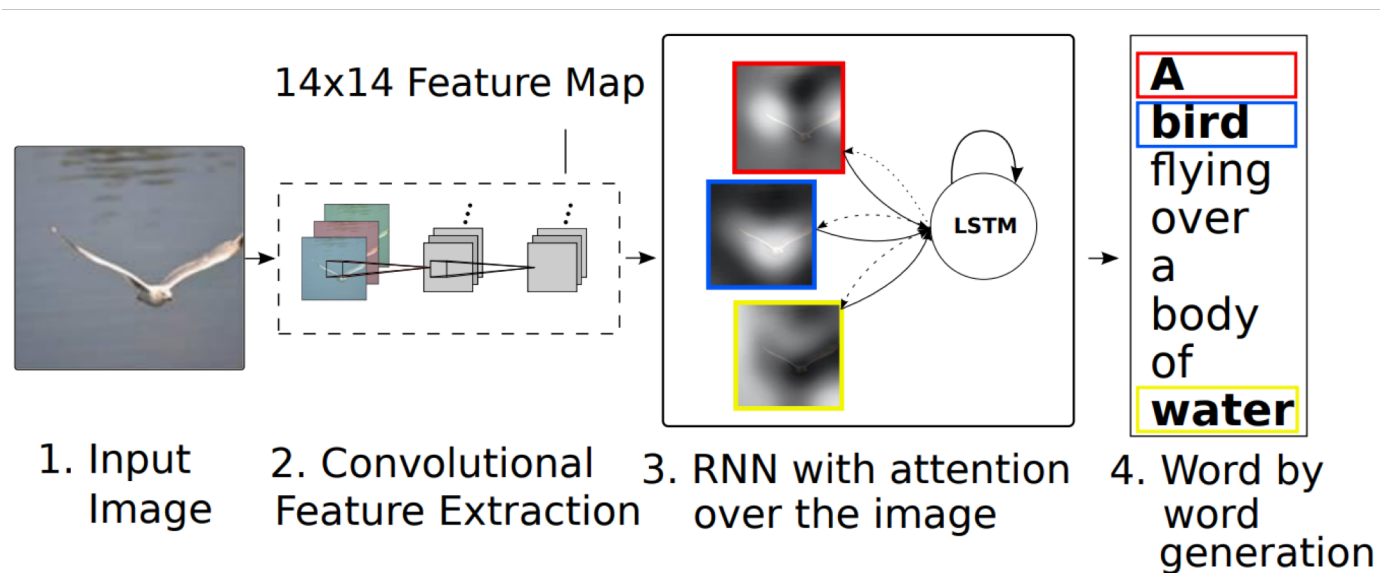
Overview Challenges Download Evaluate Leaderboard

Table-C5 Table-C40 2015 Captioning Challenge

Last update: June 8, 2015. Visit CodaLab for the latest results.

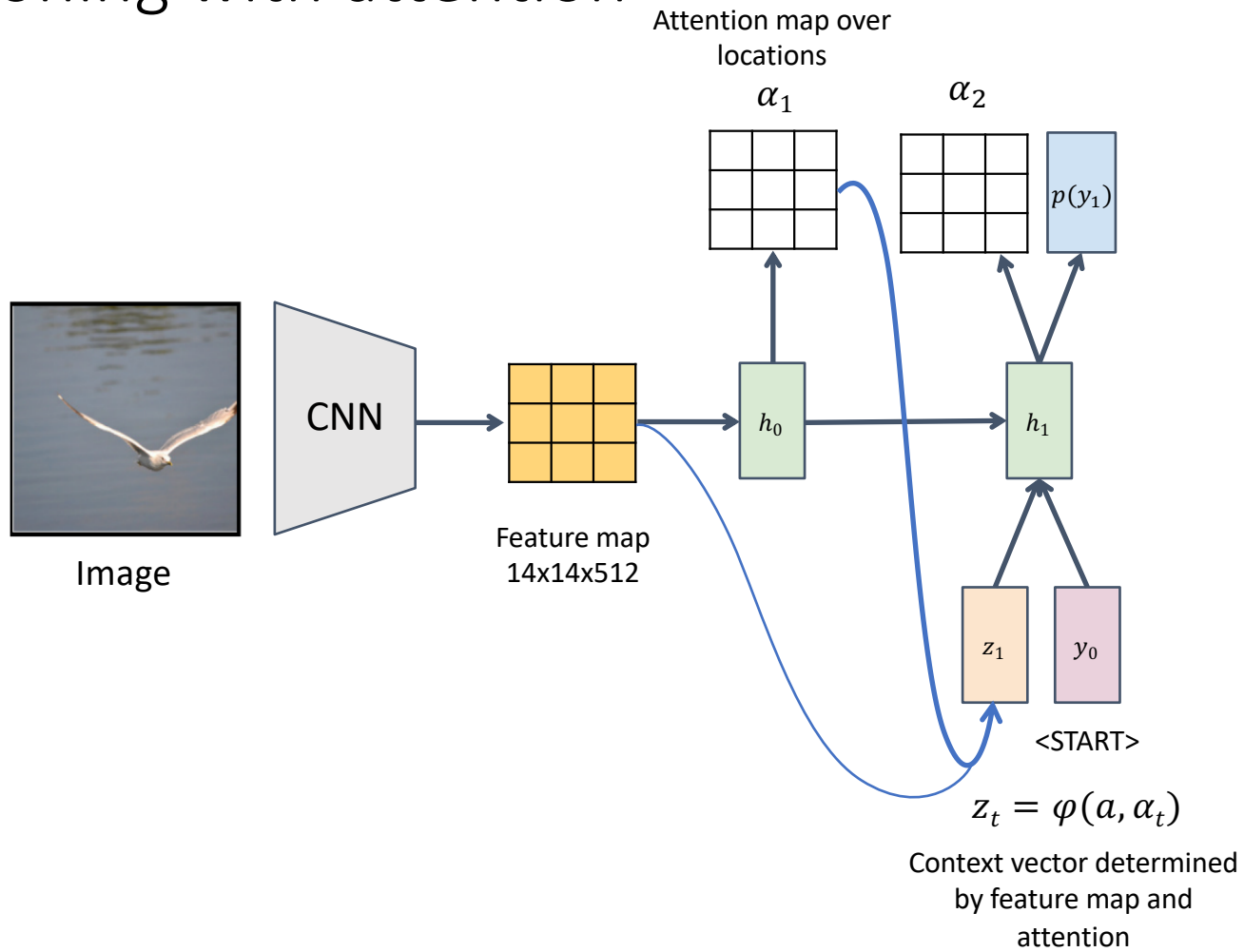
	M1	M2	M3	M4	M5
Human <sup>[5]</sup>	0.638	0.675	4.836	3.428	0.352
Google <sup>[4]</sup>	0.272	0.217	4.107	2.712	0.232
MSR <sup>[8]</sup> M1	Percentage of captions that are evaluated as better or equal to human caption.				
Montreal M2	Percentage of captions that pass the Turing Test.				
MSR Ca M3	Average correctness of the captions on a scale 1-5 (incorrect - correct).				
Berkeley M4	Average amount of detail of the captions on a scale 1-5 (lack of details - very detailed).				
m-RNN <sup>[1]</sup> M5	Percentage of captions that are similar to human description.				
Nearest Neighbor <sup>[11]</sup>	0.216	0.255	3.801	2.716	0.196
PicSOM <sup>[13]</sup>	0.202	0.250	3.965	2.552	0.182
Brno University <sup>[3]</sup>	0.194	0.213	3.079	3.482	0.154
m-RNN (Baidu/ UCLA) <sup>[16]</sup>	0.190	0.241	3.831	2.548	0.195
MIL <sup>[6]</sup>	0.168	0.197	3.349	2.915	0.159
MLBL <sup>[7]</sup>	0.167	0.196	3.659	2.420	0.156

# Captioning with attention



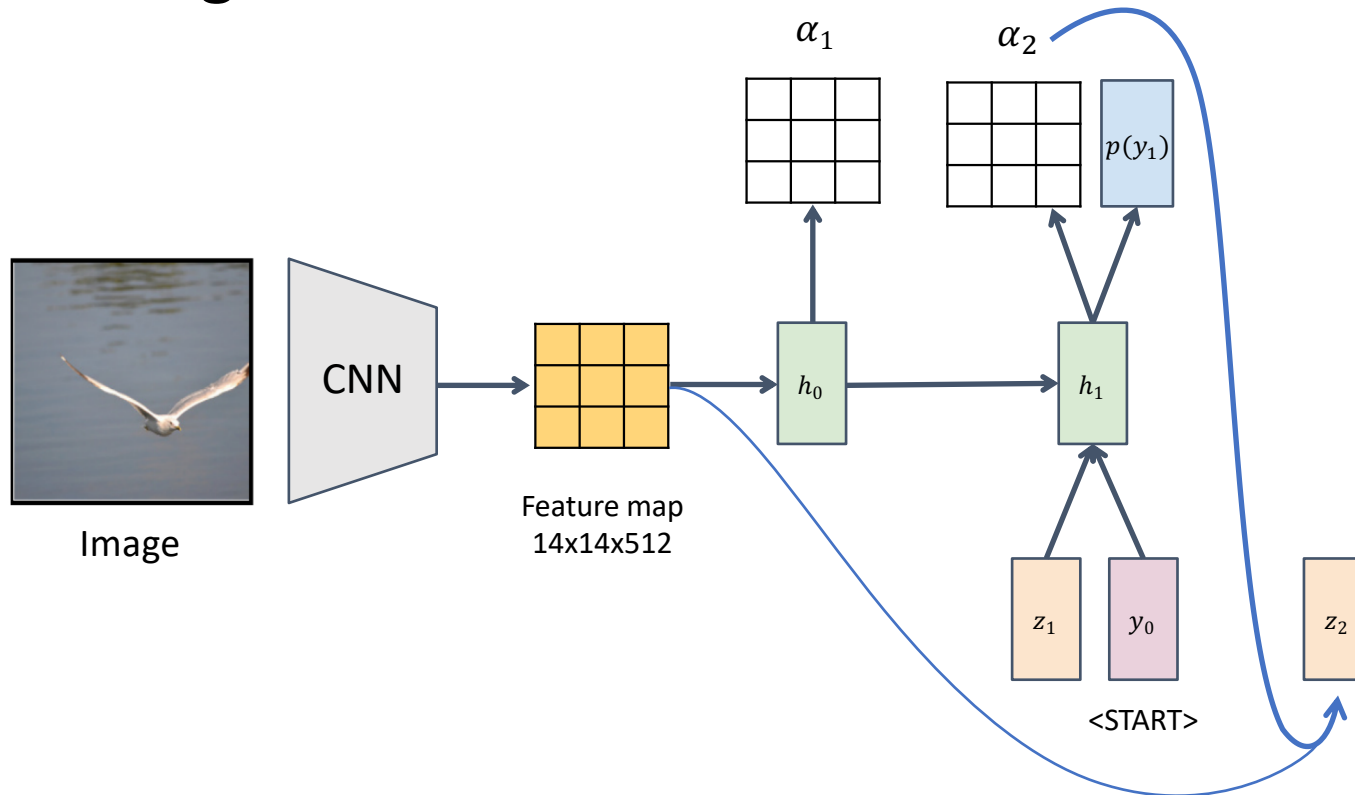
K. Xu et al., [Show, Attend and Tell: Neural Image Caption Generation with Visual Attention](#), ICML 2015

# Captioning with attention

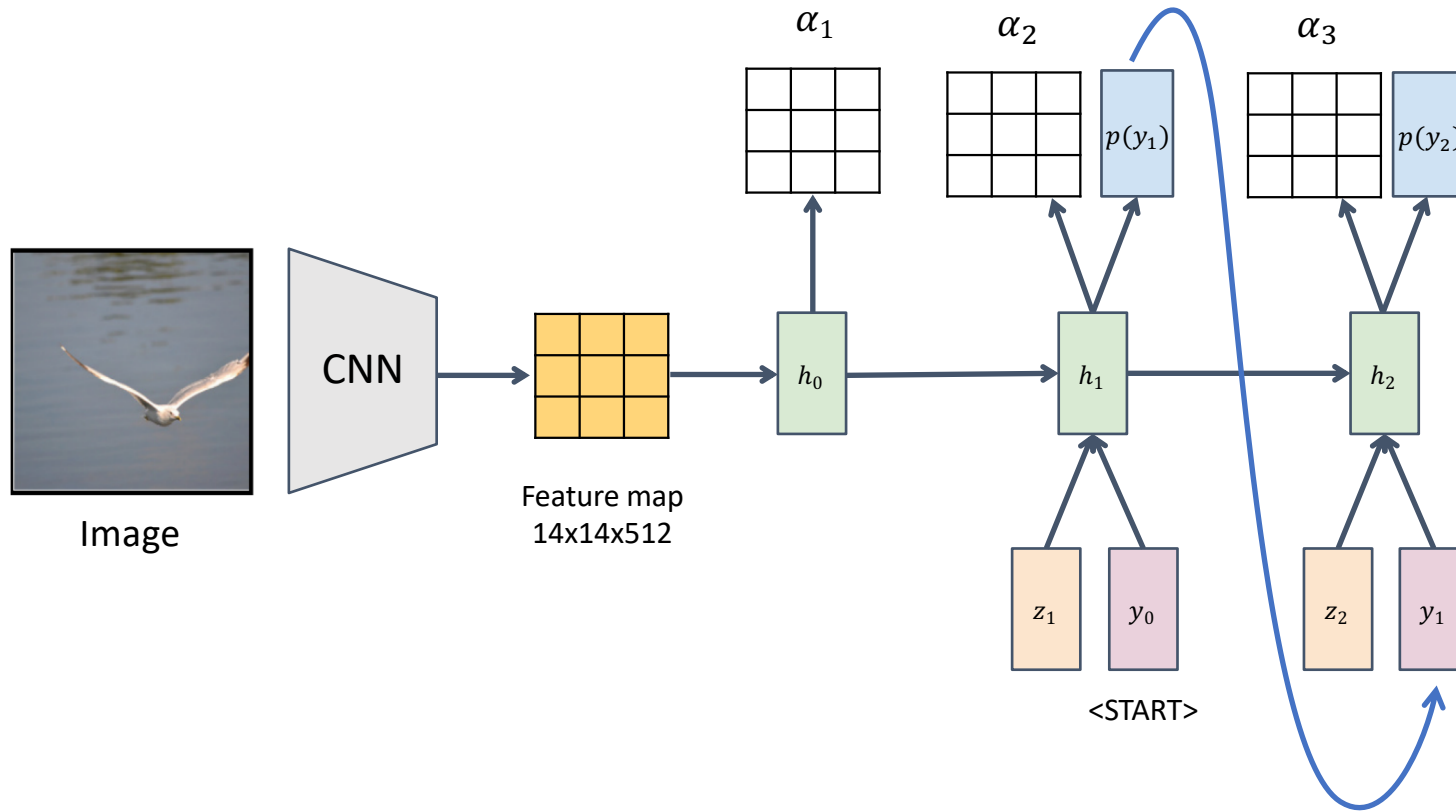


Adapted from [Stanford CS231n](#), [Berkeley CS294](#)

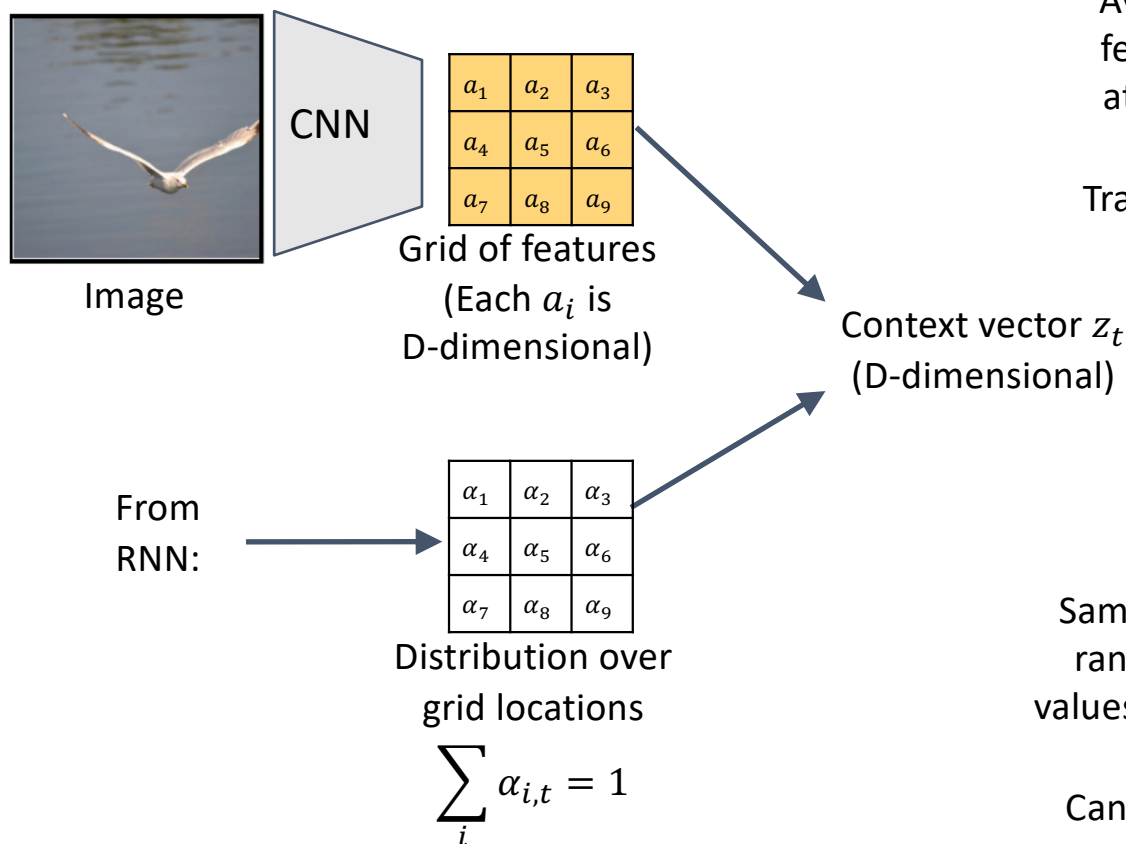
# Captioning with attention



# Captioning with attention



# “Soft” and “hard” attention



## Soft attention:

Average over locations of feature map weighted by attention:  $z_t = \sum_i \alpha_{i,t} a_i$

Train with gradient descent

## Hard attention:

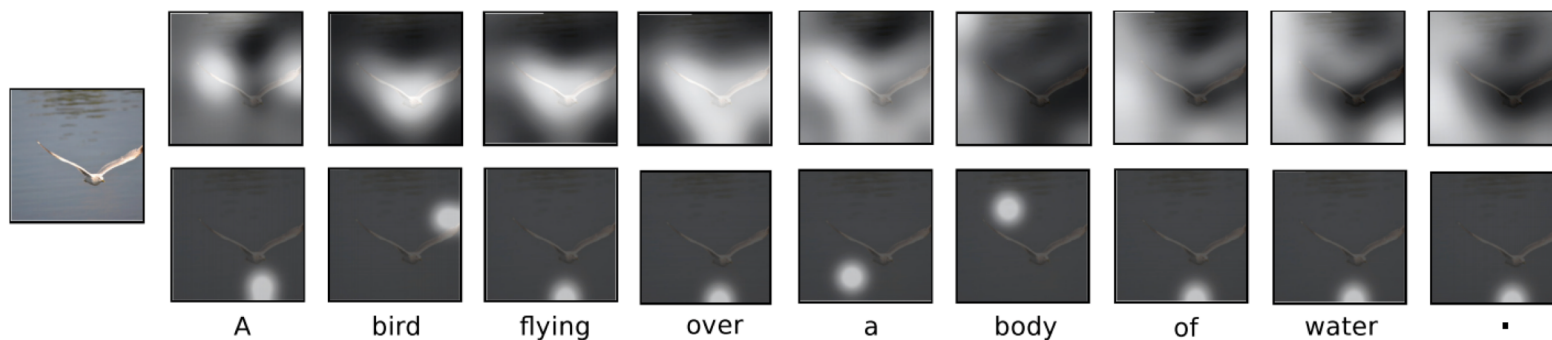
Sample ONE location:  $z_t$  is a random variable taking on values  $a_i$  with probabilities  $\alpha_{i,t}$

Can't use gradient descent; need reinforcement learning

# “Soft” and “hard” attention

## Soft attention:

Average over locations of feature map weighted by attention:  $z_t = \sum_i \alpha_{i,t} a_i$



## Hard attention:

Sample ONE location:  $z_t$  is a random variable taking on values  $a_i$  with probabilities  $\alpha_{i,t}$



# Example Results

- C



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.

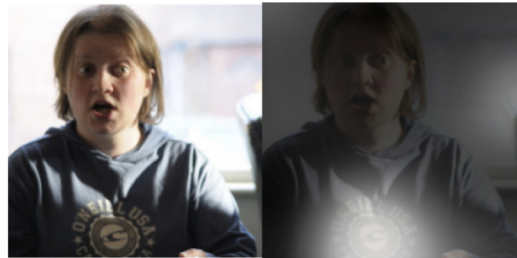
---

# Example Results

- M



A large white bird standing in a forest.



A woman holding a clock in her hand.



A man wearing a hat and a hat on a skateboard.



A person is standing on a beach with a surfboard.



A woman is sitting at a table with a large pizza.



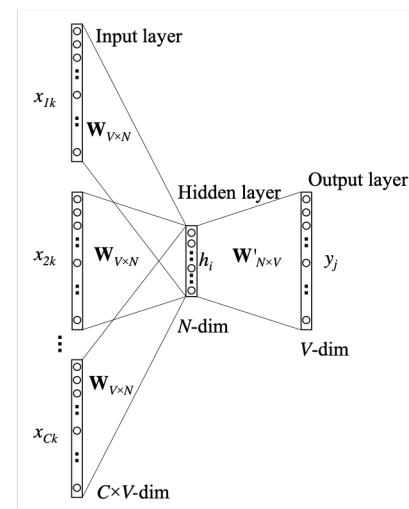
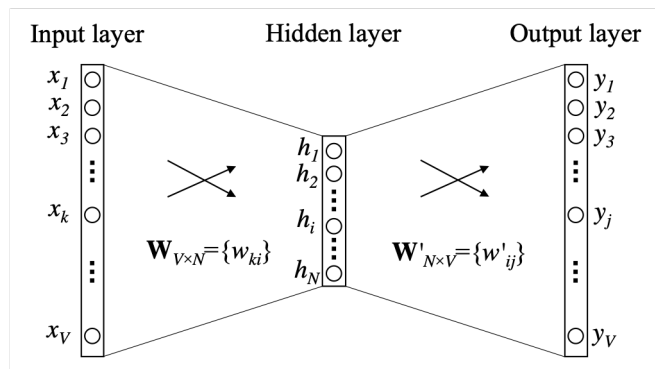
A man is talking on his cell phone while another man watches.

# Word embeddings

- Mikolov – how to learn word representations
- Skip gram model – learn quality representations of words from unstructured text
- $\text{vec}(\text{“Madrid”}) - \text{vec}(\text{“Spain”}) + \text{vec}(\text{“France”})$  is closer to  $\text{vec}(\text{“Paris”})$
- Distributed Representations of Words and Phrases and their Compositionality
- Idea – how to capture the similarity between words instead of treating them as atomic units
- Train a network to generate vector representations of words

# word2vec

- Word2vec parameter learning explained Xin Rong  
<https://arxiv.org/abs/1411.2738>
- [bit.ly/wevi-online](http://bit.ly/wevi-online)
- Variations – one word to one – hidden layer is the embedding
- N-gram model to one
- Other alternative GloVec



# RNN vs Sequence to Sequence models

- Especially when it comes to seq2seq models, is one hidden state really enough to capture global information pertaining to the translation?
- Idea – learn a context vector – for each input vector we learn a set of weights of how much the remainder of the sentence is affected by the rest of the sentence
- Sequence attention model – more detail next