



Perception

Image Captioning

Yezhou Yang, Ph.D.
Assistant Professor
Zhiyuan Fang, Teaching Assistant
Arizona State University

Problem Definition

$F($  $) = \text{"A furry, grey kitten."}$

| Image Captioning:

- Describing the image in an open form natural language sentence.

Captioning Examples



A dog is running on the grass with a ball.



A young lady leans on a table in the kitchen.

Captioning Examples

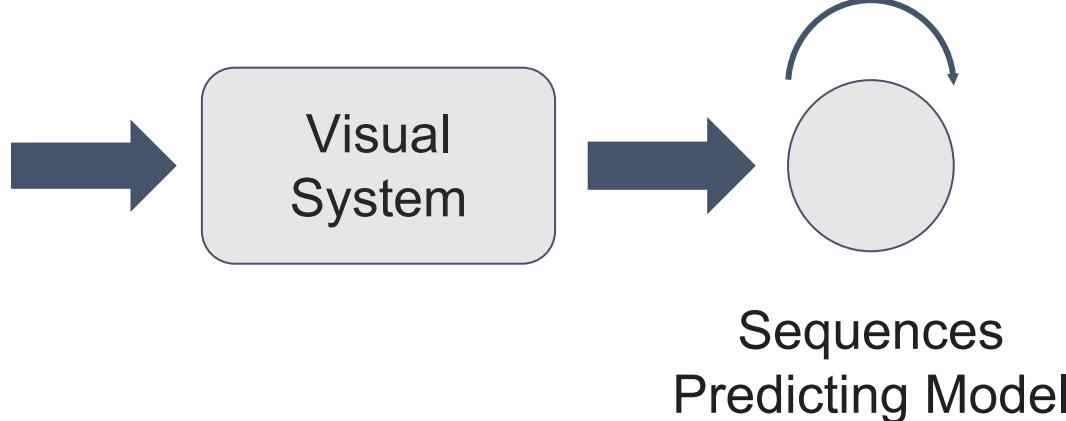


The men are playing a game of baseball on the field.



A city bus rides down a street in front of some buildings.

General Workflow



| Image Captioning:

- State-of-the-art method combines recent advances in computer vision and machine translation and that can be used to generate natural sentences describing an image.

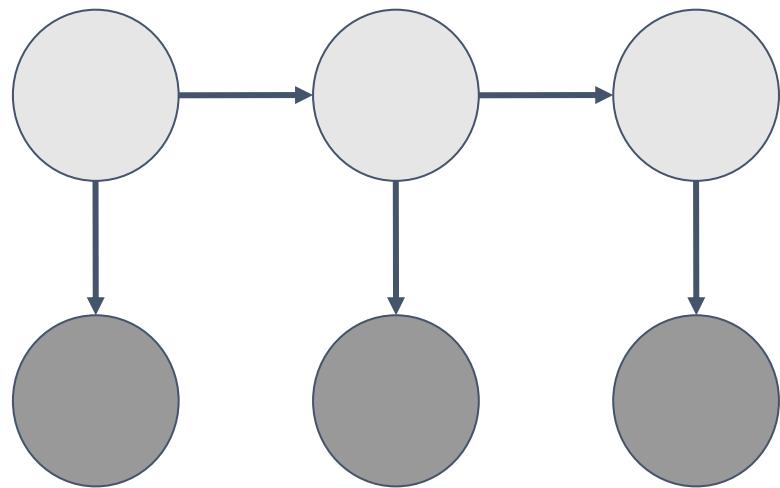
Relevant Methods

To generate the languages, we are training a sequence predicting model:

- $P(\text{next word} \mid \text{previous words})$

Hidden Markov Model (HMM):

- Is a statistical Markov model in which the system being modeled is assumed to be a Markov process with unobserved (i.e. hidden) states.



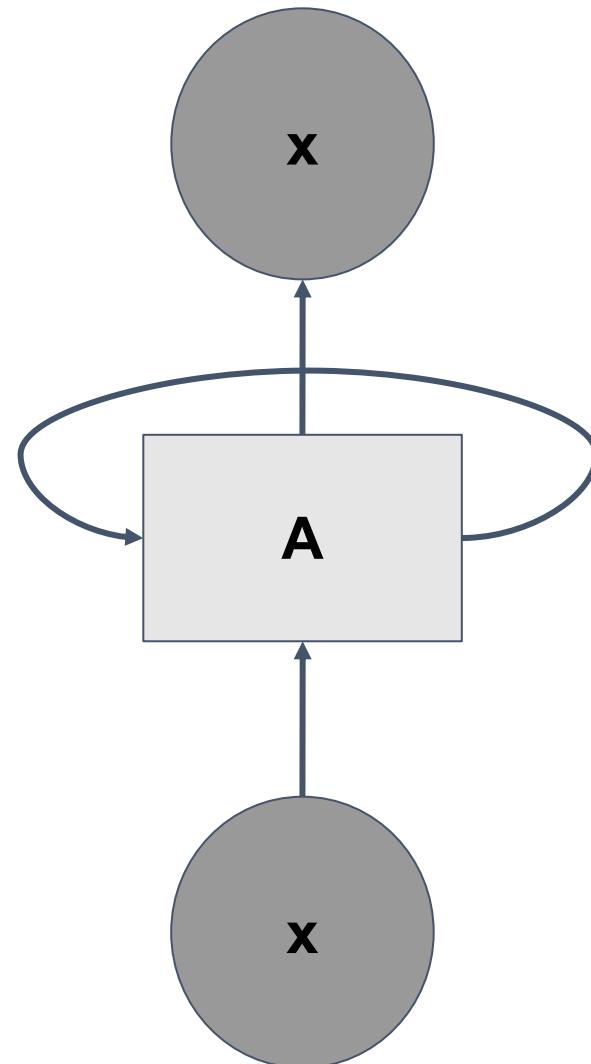
Relevant Methods

| To generate the languages, we are training a sequence predicting model:

- $P(\text{next word} \mid \text{previous words})$

| Recurrent Neural Network (RNN):

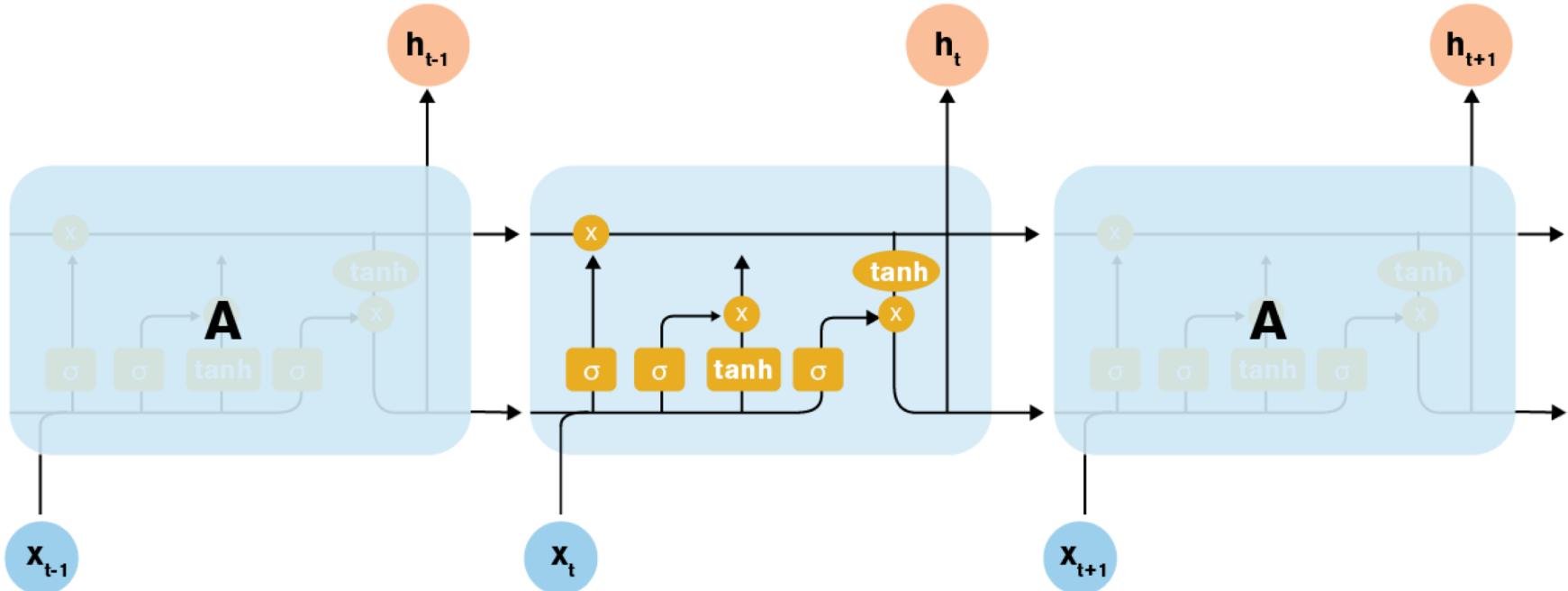
- Is a class of artificial neural network where connections between nodes form a directed graph along a temporal sequence.



Relevant Methods

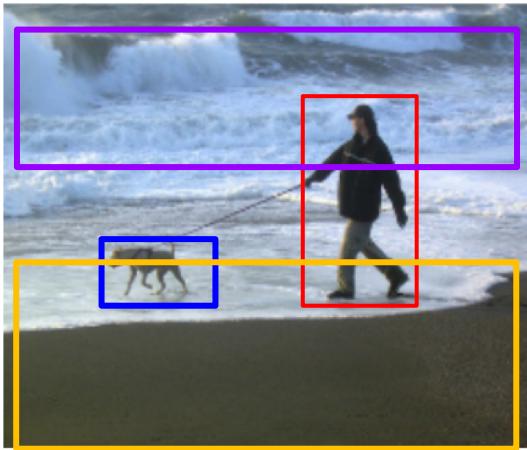
Long Short-term Memory (LSTM):

- Is a variation of RNN to deal with the exploding and **vanishing** gradient problems that can be encountered when training traditional RNNs.

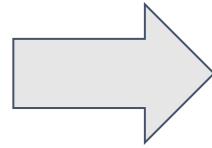


"Understanding LSTM Networks," Understanding LSTM Networks -- colah's blog. [Online]. Available: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>. [Accessed: 11-May-2019].

Relevant Methods



Visual
Classifier



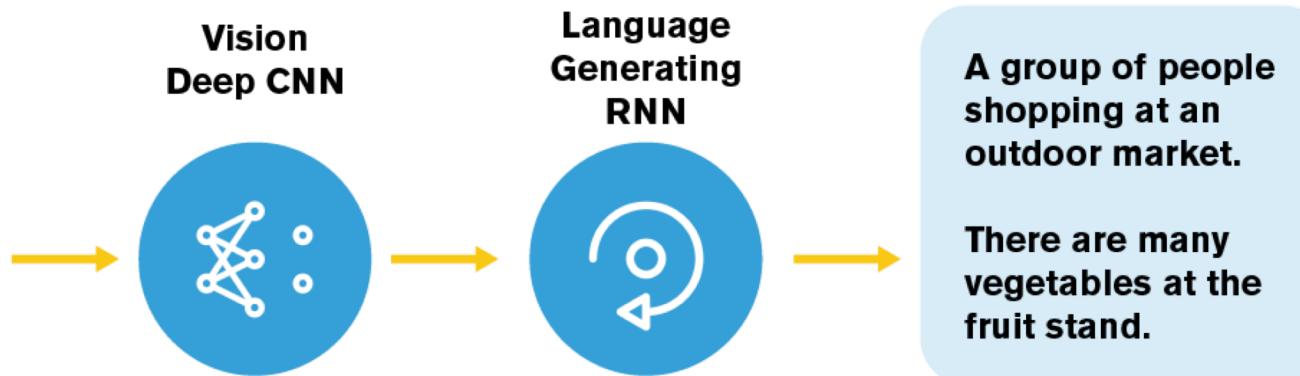
{ Person, Dog,
Coat, Sea }



HMM

"A person is walking their dog on the coast."

Relevant Methods



Reading Materials

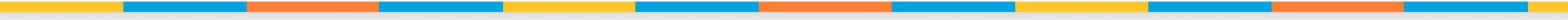


Reading Materials

- Image captioning with semantic attention
- Self-critical sequence training for image captioning
- Neural Networks and Deep Learning:Image Captioning
- Designing, Visualizing and Understanding Deep Neural Networks

You, Quanzeng, et al. "Image captioning with semantic attention." Proceedings of the IEEE conference on computer vision and pattern recognition. 2016.
Rennie, Steven J., et al. "Self-critical sequence training for image captioning." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2017.
Neural Networks and Deep Learning:Image Captioning, Mike Mozer
Designing, Visualizing and Understanding Deep Neural Networks, John Canny

Reading Materials



Reading Materials

- Automated Image Captioning with ConvNets and Recurrent Nets
- Corpus-guided sentence generation of natural images
- Show and tell: A neural image caption generator
- An Introduction to Conditional Random Fields

Automated Image Captioning with ConvNets and Recurrent Nets, Andrej Karpathy, Fei-Fei Li

Yang, Yezhou, et al. "Corpus-guided sentence generation of natural images." EMNLP, 2011.

Vinyals, Oriol, et al. "Show and tell: A neural image caption generator." Proceedings of the IEEE conference on computer vision and pattern recognition. 2015.

An Introduction to Conditional Random Fields, Charles Sutton and Andrew McCallum