

Recurrent Neural Network - PART II

Minlie Huang

aihuang@tsinghua.edu.cn
Department of Computer Science
Tsinghua University
Some slides are from Dr. Xipeng Qiu

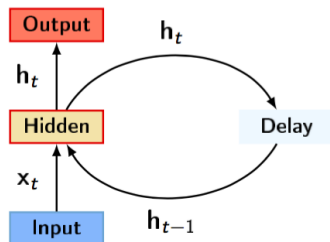
2018

Recurrent Neural Network (RNN)

The RNN has recurrent hidden states whose output at each time is dependent on that of the previous time.

Given a sequence $x(1:n) = (x_1, x_2, \dots, x_t, \dots, x_n)$ (**each x_i is a vector**), the RNN updates its recurrent hidden state $h(t)$ by

$$h_t = \begin{cases} 0 & \text{if } t=0 \\ f(h_{t-1}, x_t) & \text{otherwise} \end{cases} \quad (1)$$



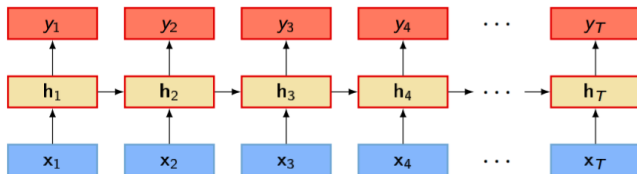
Simple Recurrent Network

- The recurrence function:

$$h_t = f(Uh_{t-1} + Wx_t + b) \quad (2)$$

where f is non-linear function (for instance, *sigmoid* or *tanh*).
Weight $U \in \mathcal{R}^{d \times d}$, weight $W \in \mathcal{R}^{d \times k}$, hidden state $h \in \mathcal{R}^d$,
and input $X \in \mathcal{R}^k$.

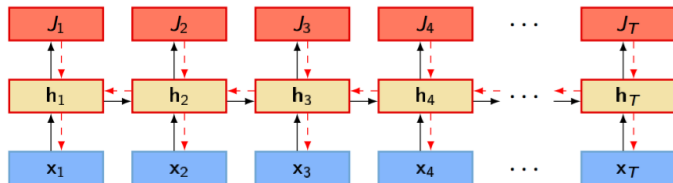
- A typical RNN structure can be illustrated as follows:



Backpropagation Through Time, BPTT

Suppose loss at time t is J_t , then the total loss is $\sum_{t=1}^T J_t$.
The gradient of J is

$$\frac{\partial J}{\partial U} = \sum_{t=1}^T \frac{\partial J_t}{\partial U} = \sum_{t=1}^T \frac{\partial h_t}{\partial U} \frac{\partial J_t}{\partial h_t} \quad (3)$$



Matrix Calculus

$$\frac{\partial y}{\partial \mathbf{x}} = \left[\frac{\partial y}{\partial x_1} \quad \frac{\partial y}{\partial x_2} \quad \cdots \quad \frac{\partial y}{\partial x_n} \right]$$

Figure: $y = f(x_1, x_2, \dots, x_n)$

$$\frac{\partial \mathbf{Y}}{\partial x} = \begin{bmatrix} \frac{\partial y_{11}}{\partial x} & \frac{\partial y_{12}}{\partial x} & \cdots & \frac{\partial y_{1n}}{\partial x} \\ \frac{\partial y_{21}}{\partial x} & \frac{\partial y_{22}}{\partial x} & \cdots & \frac{\partial y_{2n}}{\partial x} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_{m1}}{\partial x} & \frac{\partial y_{m2}}{\partial x} & \cdots & \frac{\partial y_{mn}}{\partial x} \end{bmatrix}$$

Figure: \mathbf{Y} is a matrix.

Matrix Calculus (cont'd)

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$$

Figure: \mathbf{y}, \mathbf{x} are column vectors.

$$\frac{\partial \mathbf{F}}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial \mathbf{F}}{\partial X_{1,1}} & \cdots & \frac{\partial \mathbf{F}}{\partial X_{1,n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial \mathbf{F}}{\partial X_{m,1}} & \cdots & \frac{\partial \mathbf{F}}{\partial X_{m,n}} \end{bmatrix}$$

Figure: \mathbf{F}, \mathbf{X} are matrices. 

Gradient of RNN

$$\frac{\partial J}{\partial U} = \sum_{t=1}^T \sum_{k=1}^t \frac{\partial h_k}{\partial U} \frac{\partial h_t}{\partial h_k} \frac{\partial y_t}{\partial h_t} \frac{\partial J_t}{\partial y_t} \quad (4)$$

$$\frac{\partial h_t}{\partial h_k} = \prod_{i=k+1}^t \frac{\partial h_i}{\partial h_{i-1}} = \prod_{i=k+1}^t \text{diag}[f'(h_{i-1})] * U \quad (5)$$

$$\frac{\partial J}{\partial U} = \sum_{t=1}^T \sum_{k=1}^t \frac{\partial h_k}{\partial U} \left(\prod_{i=k+1}^t \text{diag}[f'(h_{i-1})] * U \right) \frac{\partial y_t}{\partial h_t} \frac{\partial J_t}{\partial y_t} \quad (6)$$

Derive the Results

$$\vec{h}_t = f(\vec{U} * \vec{h}_{t-1} + \vec{W} * \vec{X}_t) \quad (7)$$

where weight $\vec{U} \in \mathcal{R}^{d \times d}$, weight $\vec{W} \in \mathcal{R}^{d \times k}$, hidden state $\vec{h} \in \mathcal{R}^d$, and input $\vec{X} \in \mathcal{R}^k$.

$$\begin{bmatrix} \vec{h}_t[1] \\ \vec{h}_t[2] \\ \vdots \\ \vec{h}_t[d] \end{bmatrix} = \begin{bmatrix} f(\vec{U}_{1*} \cdot \vec{h}_{t-1} + \vec{W}_{1*} \cdot \vec{X}_t) \\ f(\vec{U}_{2*} \cdot \vec{h}_{t-1} + \vec{W}_{2*} \cdot \vec{X}_t) \\ \vdots \\ f(\vec{U}_{d*} \cdot \vec{h}_{t-1} + \vec{W}_{d*} \cdot \vec{X}_t) \end{bmatrix} = \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_d \end{bmatrix} \quad (8)$$

where each f_i is a scalar since $\vec{U}_{i*} \cdot \vec{h}_{t-1} + \vec{W}_{i*} \cdot \vec{X}_t$ is a scalar.

Derive the Results (Cont'd)

$$\frac{\partial \vec{h}_t}{\partial \vec{h}_{t-1}} = \begin{bmatrix} \frac{\partial f_1}{\partial \vec{h}_{t-1}} \\ \frac{\partial f_2}{\partial \vec{h}_{t-1}} \\ \vdots \\ \frac{\partial f_d}{\partial \vec{h}_{t-1}} \end{bmatrix} = \begin{bmatrix} \vec{U}_{1*} \cdot f'(\vec{h}_{t-1})[1] \\ \vec{U}_{2*} \cdot f'(\vec{h}_{t-1})[2] \\ \vdots \\ \vec{U}_{d*} \cdot f'(\vec{h}_{t-1})[d] \end{bmatrix} \quad (9)$$

where each $f'(\vec{h}_{t-1})[i]$ is a scalar, $f'(\vec{h}_{t-1})$ is a column vector, and \vec{U}_{i*} is a row vector.

$$\frac{\partial \vec{h}_t}{\partial \vec{h}_{t-1}} = \begin{bmatrix} f'(\vec{h}_{t-1})[1] & & & \\ & f'(\vec{h}_{t-1})[2] & & \\ & & \ddots & \\ & & & f'(\vec{h}_{t-1})[d] \end{bmatrix} \cdot \begin{bmatrix} \vec{U}_{1*} \\ \vec{U}_{2*} \\ \vdots \\ \vec{U}_{d*} \end{bmatrix} \quad (10)$$

$$= \text{diag}[f'(\vec{h}_{t-1})] \cdot \vec{U} \quad (11)$$

Connections to MLP?

Why is RNN a deep structure?

Connections to MLP?

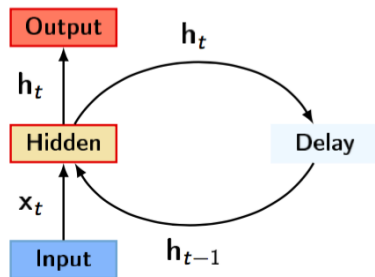


Figure: Traditional View.

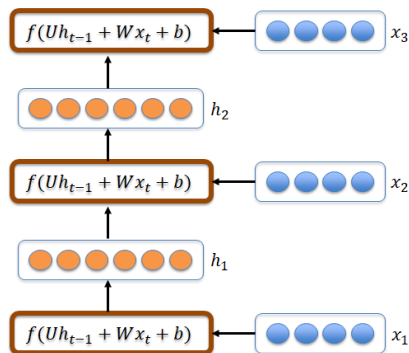


Figure: Deep View.

Long-Term Dependencies

Define $\gamma = \|\text{diag}[f'(h_{i-1})] * U\|$,

- ▶ Exploding Gradient Problem: When $\gamma > 1$ and $t - k \rightarrow \infty$, $\gamma^{t-k} \rightarrow \infty$.
- ▶ Vanishing Gradient Problem: When $\gamma < 1$ and $t - k \rightarrow \infty$, $\gamma^{t-k} \rightarrow 0$.

There are various ways to solve Long-Term Dependency problem.

Gradient Explosion & Vanishing

- ▶ Gradient norm clipping (Mikolov thesis 2012; Pascanu, Mikolov, Bengio, ICML 2013)
-

$$\hat{g} \leftarrow \frac{\partial \mathcal{J}}{\partial \theta}$$

IF

$$\|\hat{g}\| > \textit{Threshold}$$

THEN

$$\hat{g} \leftarrow \frac{\textit{Threshold}}{\|\hat{g}\|} * \hat{g}$$

-
- ▶ Gradient propagation regularizer (avoid vanishing gradient)
 - ▶ LSTM self-loops (avoid vanishing gradient)

Long Short-Term Memory (LSTM)

Key difference to RNN: a memory cell c which is controlled by three gates:

- ▶ input gate i :

$$i_t = \sigma(W_i x_t + U_i h_{t-1} + b_i) \quad (12)$$

- ▶ output gate o :

$$o_t = \sigma(W_o x_t + U_o h_{t-1} + b_o) \quad (13)$$

- ▶ forget gate f :

$$f_t = \sigma(W_f x_t + U_f h_{t-1} + b_f) \quad (14)$$

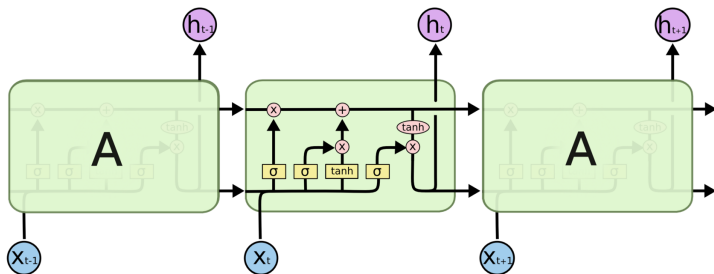
- ▶ update of memory cell and hidden state:

$$\widetilde{C}_t = \tanh(W_c x_t + U_c h_{t-1} + b_c) \quad (15)$$

$$C_t = f_t \otimes C_{t-1} + i_t \otimes \widetilde{C}_t \quad (16)$$

$$h_t = o_t \otimes \tanh(C_t) \quad (17)$$

LSTM Architecture



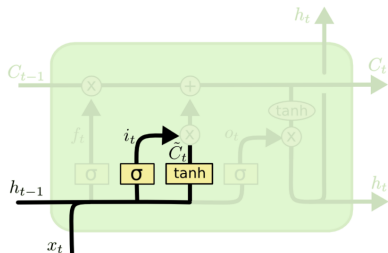
$$\widetilde{C}_t = \tanh(W_c x_t + U_c h_{t-1} + b_c)$$

$$C_t = f_t \otimes C_{t-1} + i_t \otimes \widetilde{C}_t$$

$$h_t = o_t \otimes \tanh(C_t)$$

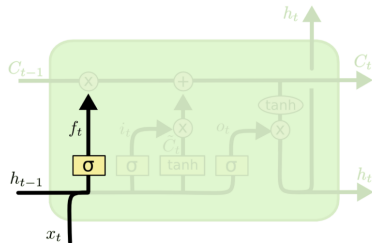
The picture and following 4 pictures are from
<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>.

LSTM: Input and Forget Gates



$$i_t = \sigma(W_i x_t + U_i h_{t-1} + b_i)$$

$$\tilde{C}_t = \tanh(W_c x_t + U_c h_{t-1} + b_c)$$

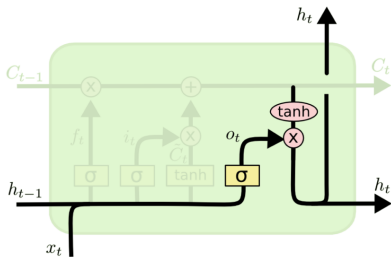


$$f_t = \sigma(W_f x_t + U_f h_{t-1} + b_f)$$

The picture is from

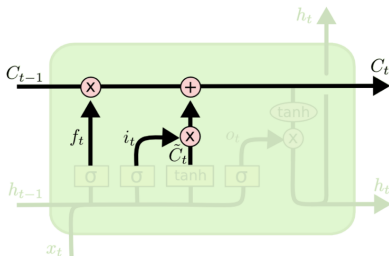
<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>.

LSTM: Output Gates and State Update



$$o_t = \sigma(W_o x_t + U_o h_{t-1} + b_o)$$

$$h_t = o_t \otimes \tanh(C_t)$$



$$C_t = f_t \otimes C_{t-1} + i_t \otimes \tilde{C}_t$$

The picture is from

<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Gated Recurrent Unit, GRU

Instead of using three gates, only two gates is employed:

- update gate z :

$$z_t = \sigma(W_z x_t + U_z h_{t-1}) \quad (18)$$

- reset gate r :

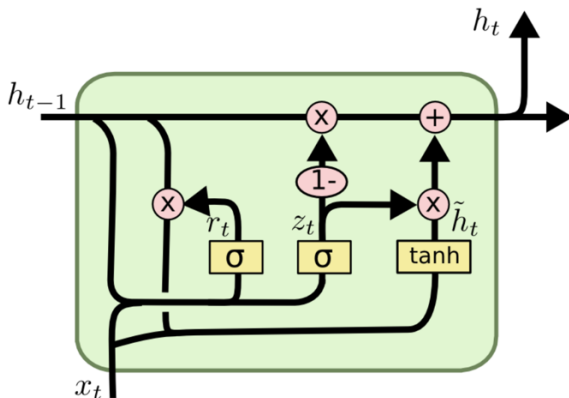
$$r_t = \sigma(W_r x_t + U_r h_{t-1}) \quad (19)$$

- state update:

$$\tilde{h}_t = \tanh(W_c x_t + U(r_t \otimes h_{t-1})) \quad (20)$$

$$h_t = (1 - z_t) \otimes h_{t-1} + z_t \otimes \tilde{h}_t \quad (21)$$

GRU Architecture



The picture is from

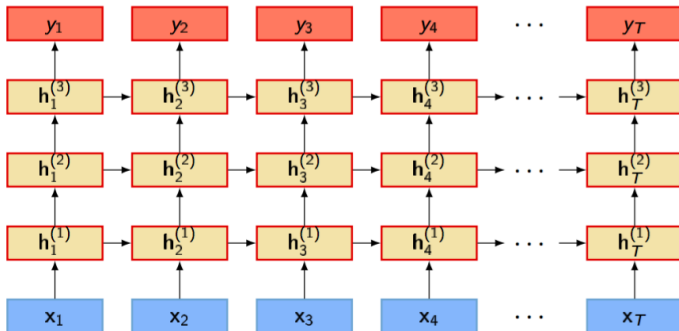
<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>.

More Exploring to LSTM...

- ▶ **How does LSTM deal with the explosion and vanishing gradient problems?**
- ▶ Refer to: Hochreiter S, Schmidhuber J. Long short-term memory[J]. Neural computation, 1997, 9(8): 1735-1780.
- ▶ **Bonus Homework (+5):** Read the paper and write a Chinese report to answer the question.

Stacked RNN/LSTM

- Stack multiple RNN/LSTM in vertical direction



- More common is seen that raw input (x_i) is connected to all layers as input.

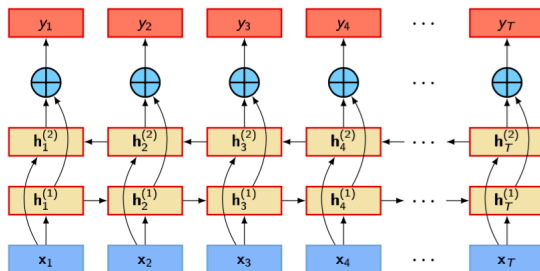
Bidirectional RNN/LSTM

- Modeling the prefix context and suffix context at the same time.

$$h_t^{(1)} = f(U^{(1)}h_{t-1}^{(1)} + W^{(1)}x_t + b^{(1)}) \quad (22)$$

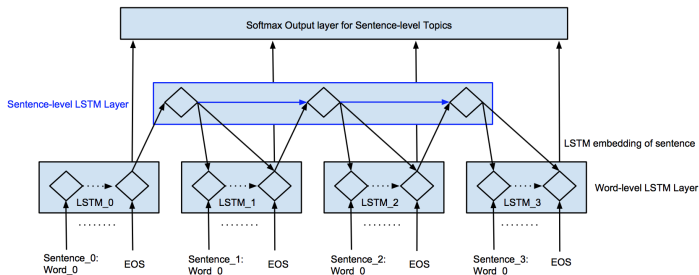
$$h_t^{(2)} = f(U^{(2)}h_{t+1}^{(2)} + W^{(2)}x_t + b^{(2)}) \quad (23)$$

$$h_t = h_t^{(1)} \oplus h_t^{(2)} \quad (24)$$



Hierarchical LSTM

Hierarchical LSTM (low-level LSTMs are taken as input to the high-level LSTM)

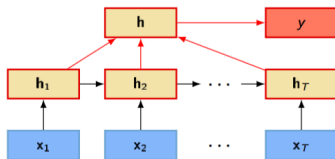


Ghosh S, Vinyals O, Strobe B, et al. Contextual LSTM (CLSTM) models for Large scale NLP tasks[J]. arXiv preprint arXiv:1602.06291, 2016.

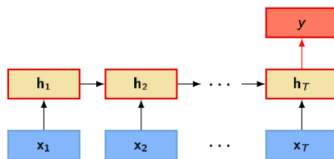
Application of RNN: Classification

Applicable to many classification problems.

- ▶ Text Classification
- ▶ Sentiment Classification



(c) Mean

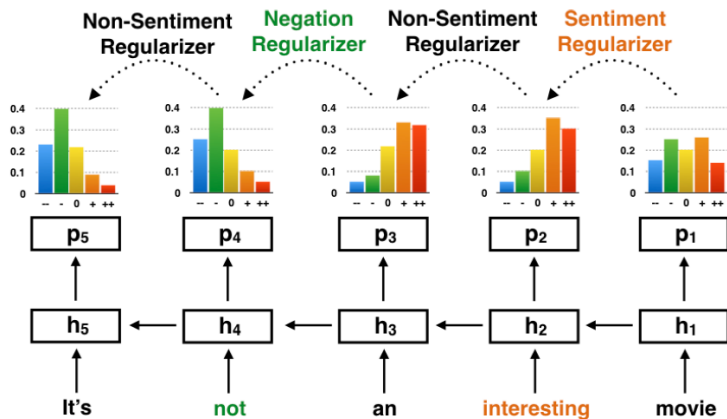


(d) Last

The left: can be enhanced with attention mechanisms.

Application of RNN: Classification

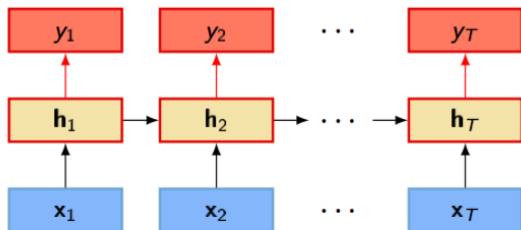
► Example: Sentiment Classification



Qiao Qian, Minlie Huang, Jinhao Lei, Xiaoyan Zhu. Linguistically Regularized LSTM for Sentiment Classification. ACL 2017.

Application of RNN: Sequence Labeling

Sequence Labeling, such as Chinese word segmentation, part-of-speech tagging, semantic role labeling

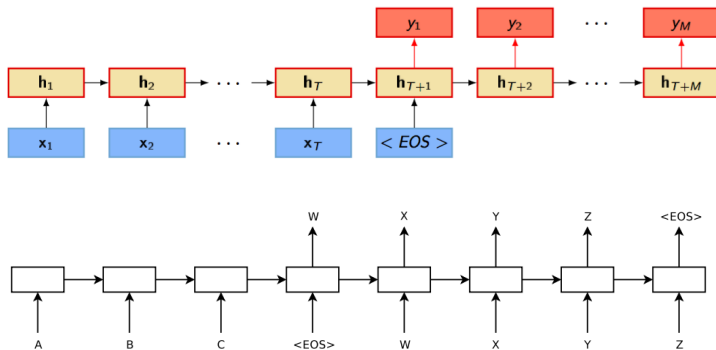


Linguistic Tools:

- ▶ Natural Language Toolkit (NLTK)
- ▶ Stanford CoreNLP
- ▶ THU Lexical Analyzer for Chinese (THULAC)

Application of RNN: Sequence to Sequence Learning

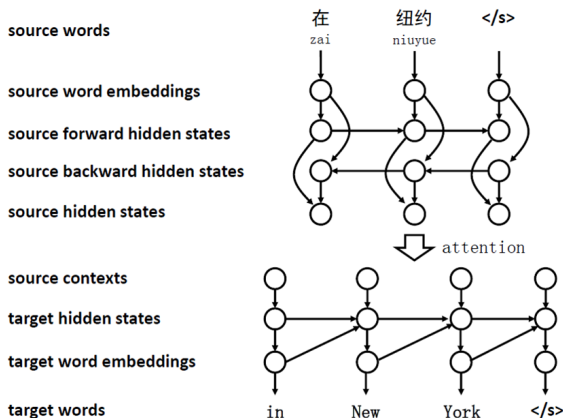
- ▶ Machine Translation
- ▶ Sequence to sequence generation (for dialogue, question answering, etc.)
- ▶ Image captioning: from image to text



Sutskever I, Vinyals O, Le Q V. Sequence to sequence learning with neural networks[C]//NIPS. 2014: 3104-3112.

Application of RNN: Sequence to Sequence Learning

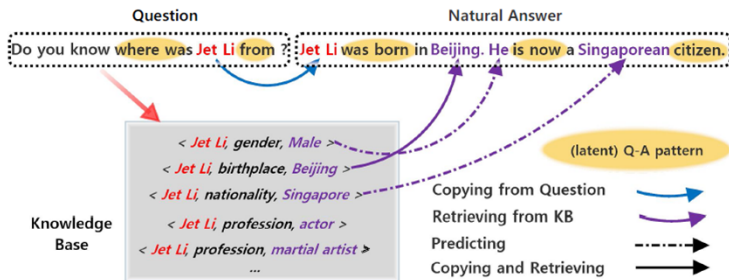
► Example: Machine Translation



Yanzhuo Ding, Yang Liu, Huanbo Luan, Maosong Sun. Visualizing and Understanding Neural Machine Translation. ACL 2017.

Application of RNN: Sequence to Sequence Learning

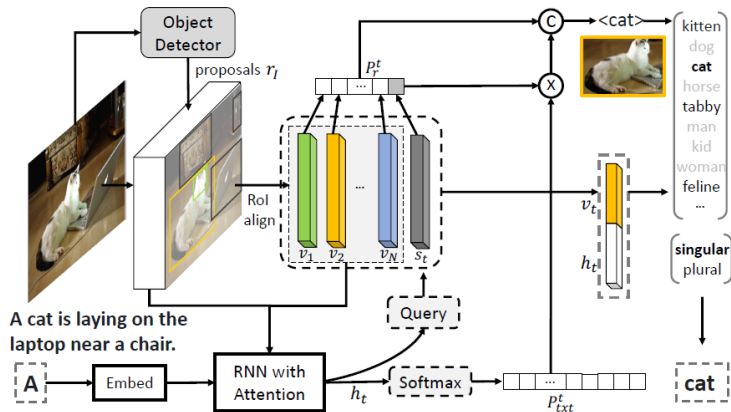
► Example: Question Answering



Shizhu He, Cao Liu, Kang Liu and Jun Zhao. Generating Natural Answers by Incorporating Copying and Retrieving Mechanisms in Sequence-to-Sequence Learning. ACL 2017.

Application of RNN: Sequence to Sequence Learning

► Example: Image Caption



Jiasen Lu, Jianwei Yang, Dhruv Batra, Devi Parikh. Neural Baby Talk. CVPR 2018.

Summary

- ▶ RNN is very good at modeling text/language
- ▶ Basic RNN models: RNN, GRU, LSTM
 - ▶ Stacked RNN
 - ▶ Bi-directional RNN
 - ▶ Hierarchical RNN
- ▶ Typical Learning Problems with RNN
 - ▶ Classification
 - ▶ Sequential Labeling
 - ▶ Sequence to Sequence Learning: MT, Dialogue Generation