

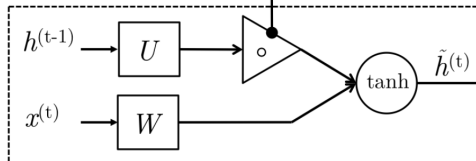
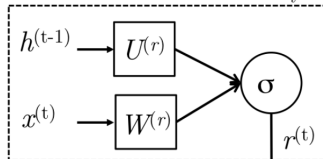
Deep Learning Technology and Application

Ge Li

Peking University

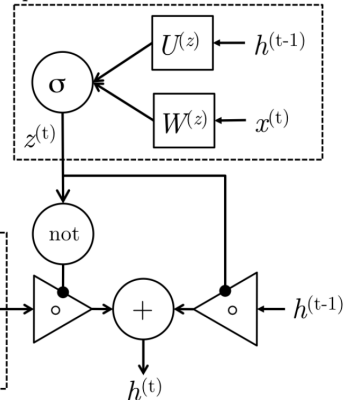
循环神经网络的发展

Reset: Include $h^{(t-1)}$ in new memory?

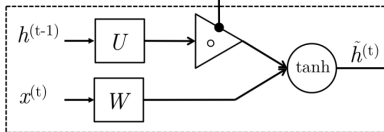
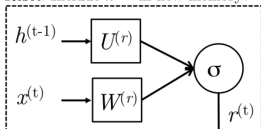


New memory: Compute new memory based on current word input $x^{(t)}$ and potentially $h^{(t-1)}$

Update: How much $h^{(t-1)}$ in next state?

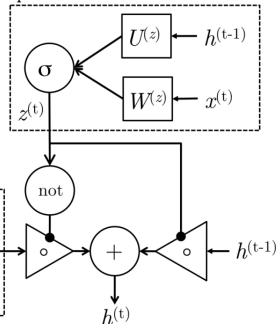


Reset: Include $h^{(t-1)}$ in new memory?



New memory: Compute new memory based on current word input $x^{(t)}$ and potentially $h^{(t-1)}$

Update: How much $h^{(t-1)}$ in next state?



$$z^{(t)} = \sigma(W^{(z)}x^{(t)} + U^{(z)}h^{(t-1)})$$

(Update gate)

$$r^{(t)} = \sigma(W^{(r)}x^{(t)} + U^{(r)}h^{(t-1)})$$

(Reset gate)

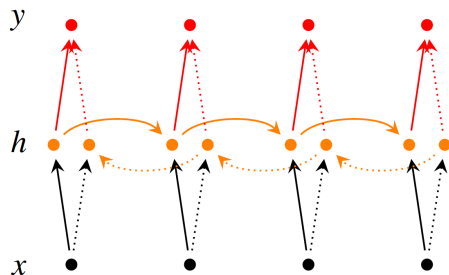
$$\tilde{h}^{(t)} = \tanh(r^{(t)} \circ U h^{(t-1)} + W x^{(t)})$$

(New memory)

$$h^{(t)} = (1 - z^{(t)}) \circ \tilde{h}^{(t)} + z^{(t)} \circ h^{(t-1)}$$

(Hidden state)

Bi-directional RNN

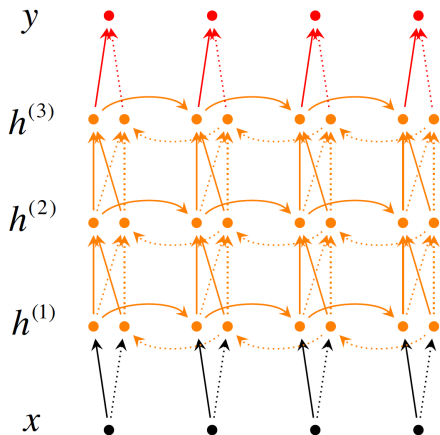


$$\vec{h}_t = f(\vec{W}x_t + \vec{V}\vec{h}_{t-1} + \vec{b})$$

$$\overleftarrow{h}_t = f(\overleftarrow{W}x_t + \overleftarrow{V}\overleftarrow{h}_{t+1} + \overleftarrow{b})$$

$$\hat{y}_t = g(Uh_t + c) = g(U[\vec{h}_t; \overleftarrow{h}_t] + c)$$

Bi-directional RNN



$$\vec{h}_t^{(i)} = f(\vec{W}^{(i)} h_t^{(i-1)} + \vec{V}^{(i)} \vec{h}_{t-1}^{(i)} + \vec{b}^{(i)})$$

$$\overleftarrow{h}_t^{(i)} = f(\overleftarrow{W}^{(i)} h_t^{(i-1)} + \overleftarrow{V}^{(i)} \overleftarrow{h}_{t+1}^{(i)} + \overleftarrow{b}^{(i)})$$

$$\hat{y}_t = g(Uh_t + c) = g(U[\vec{h}_t^{(L)}; \overleftarrow{h}_t^{(L)}] + c)$$

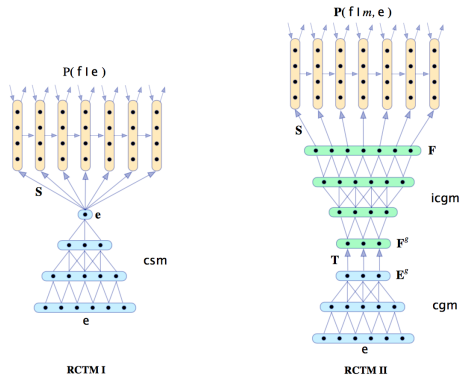
Encoder-Decoder Roadmap - 2013

[PDF] Recurrent Continuous Translation Models.

[N Kalchbrenner](#), [P Blunsom](#) - EMNLP, 2013 - [anthology.aclweb.org](#)

Abstract We introduce a class of probabilistic continuous translation models called Recurrent Continuous Translation Models that are purely based on continuous representations for words, phrases and sentences and do not rely on alignments or phrasal translation units. The models have a generation and a conditioning aspect. The generation of the translation is modelled with a target Recurrent Language Model, whereas the ...

[Cited by 279](#) [Related articles](#) [All 7 versions](#) [Cite](#) [Save](#) [More](#)



Encoder-Decoder Roadmap - 2014 RNNenc

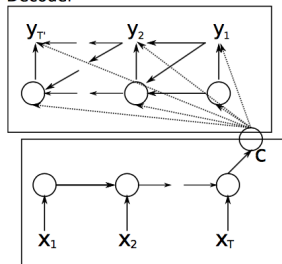
Learning phrase representations using RNN encoder-decoder for statistical machine translation

[K Cho](#), [B Van Merriënboer](#), [C Gulcehre](#)... - arXiv preprint arXiv: ..., 2014 - arxiv.org

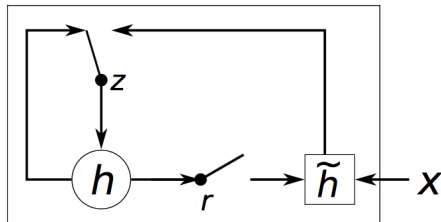
Abstract: In this paper, we propose a novel neural network model called RNN Encoder-Decoder that consists of two recurrent neural networks (RNN). One RNN encodes a sequence of symbols into a fixed-length vector representation, and the other decodes the representation into another sequence of symbols. The encoder and decoder of the proposed model are jointly trained to maximize the conditional probability of a target sequence ...

[Cited by 981](#) [Related articles](#) [All 19 versions](#) [Cite](#) [Save](#)

Decoder



Encoder



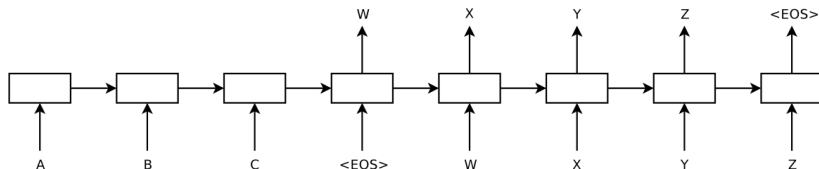
Encoder-Decoder Roadmap - 2014 NIPS

[PDF] Sequence to sequence learning with neural networks

I Sutskever, O Vinyals, QV Le - Advances in neural information ..., 2014 - papers.nips.cc

Page 1. **Sequence to Sequence Learning** with Neural Networks Ilya Sutskever Google
ilyasu@google.com ... In this paper, we present a general end-to-end approach to **sequence learning** that makes minimal assumptions on the **sequence** structure. ...

Cited by 1505 Related articles All 15 versions Cite Save More

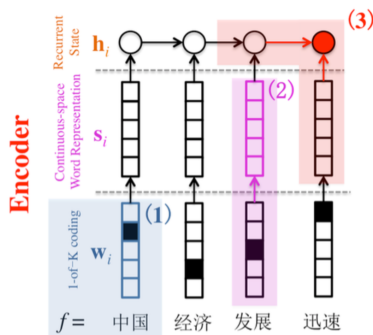
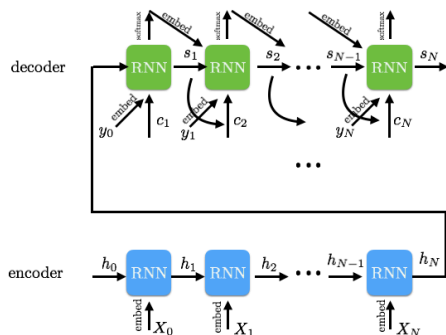


Encoder-Decoder

Encoder:

- 设 $D = (x^1, y^1), \dots, (x^N, y^N)$ 为包含 N 个平行句子的平行语料库；
(下面先针对一组平行句子进行讨论，此时，可以省去上标 N)
- 设 h_t 为 Encoding 过程中 t 时刻隐藏层的状态；

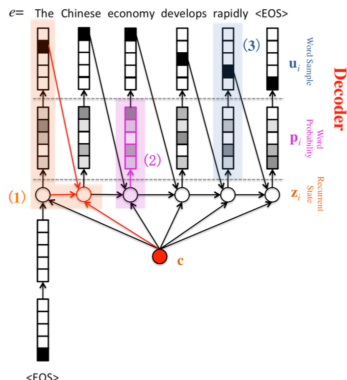
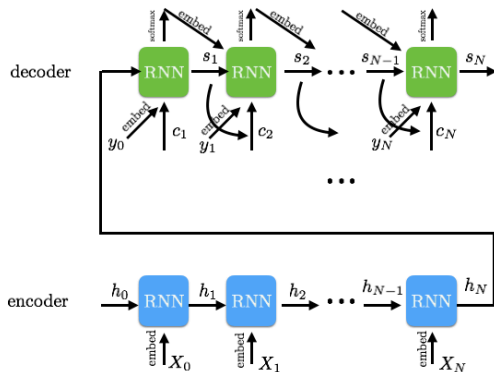
$$h_i = f(h_{i-1}, x_i)$$



Encoder-Decoder

Decoder:

- 设 h_t 为 Encoding 过程中 t 时刻隐藏层的状态；
- 设 s_o 为 Decoding 过程中 o 时刻隐藏层的状态；
- 设 c_o 为 Decoding 过程中 o 时刻的上下文信息；



Encoder-Decoder

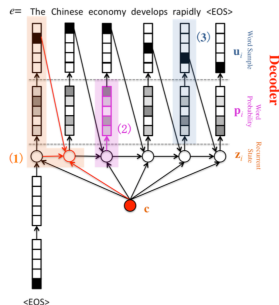
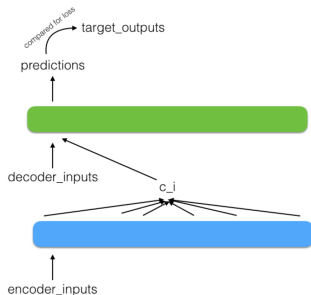
Decoder:

$$p(y_1, \dots, y_O | x_1, \dots, x_T) = \prod_{o=1}^O p(y_o | y_1, \dots, y_{o-1}, c)$$

$$p(y_o | y_1, \dots, y_{o-1}, c) = g(y_{o-1}, s_o, c)$$

$$s_o = f(y_{o-1}, s_{o-1}, c)$$

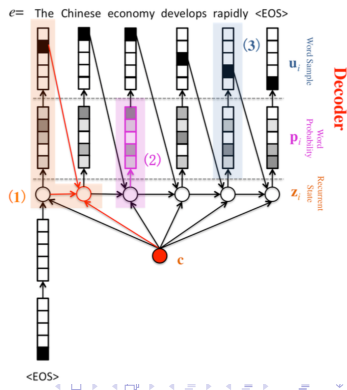
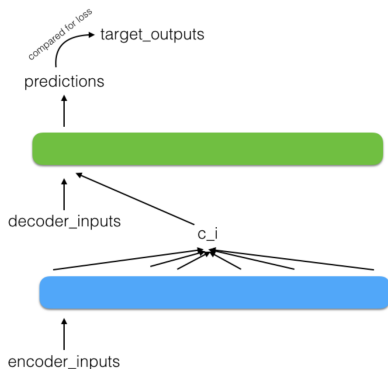
$$c = q(\{h_0, \dots, h_T\}) \quad \text{不妨先设: } c_t = h_T$$



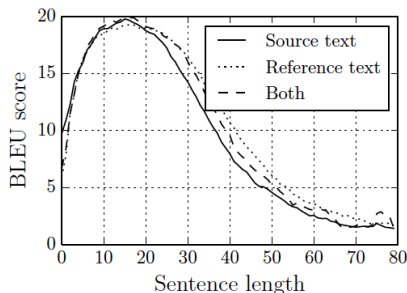
Encoder-Decoder

Decoder: 对全部语料库 $D = (x^1, y^1), \dots, (x^N, y^N)$, 训练目标为：

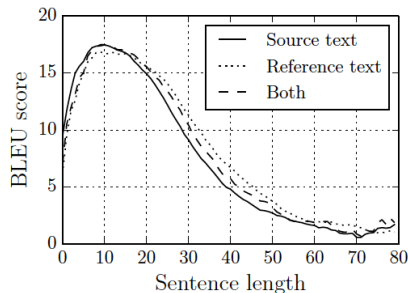
$$J(D, \Theta) = \frac{1}{N} \sum_{n=1}^N \log p(y^n | x^n, \Theta) = \frac{1}{N} \sum_{n=1}^N \sum_{o=1}^O \log p(y_o^n | y_1^n, \dots, y_{o-1}^n, c, \Theta)$$



Problem of Encoder-Decoder

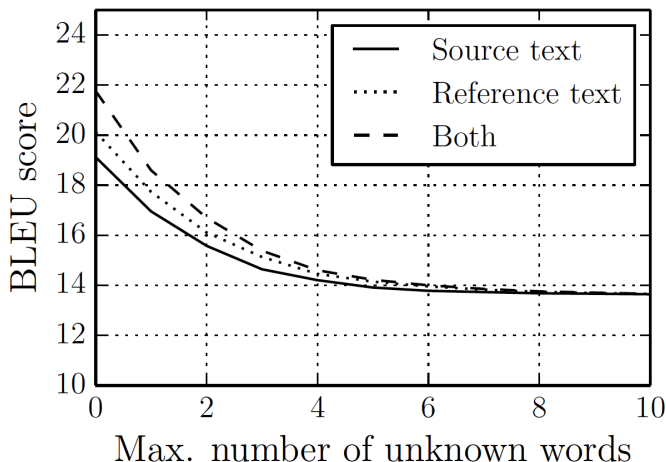


(a) RNNenc



(b) grConv

Problem of Encoder-Decoder



Attention Roadmap - 2014 RNNsearch

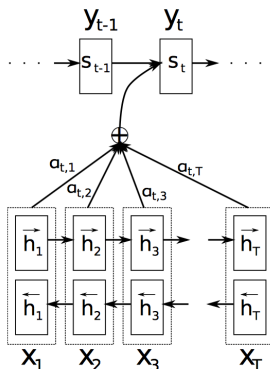
Bidirectional RNN for Annotating Sequence

Neural machine translation by jointly learning to align and translate

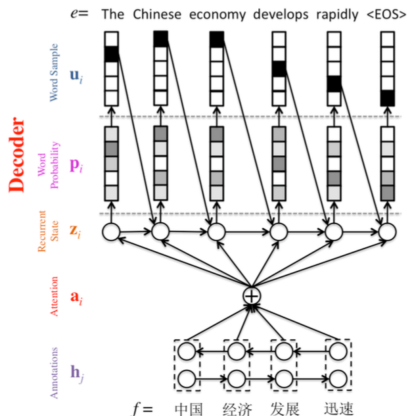
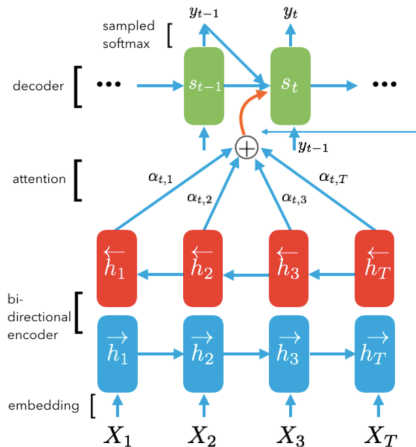
[D Bahdanau](#), [K Cho](#), [Y Bengio](#) - arXiv preprint arXiv:1409.0473, 2014 - arxiv.org

Abstract: **Neural machine translation** is a recently proposed approach to **machine translation**. Unlike the traditional statistical **machine translation**, the **neural machine translation** aims at building a single **neural** network that can be **jointly** tuned to maximize the

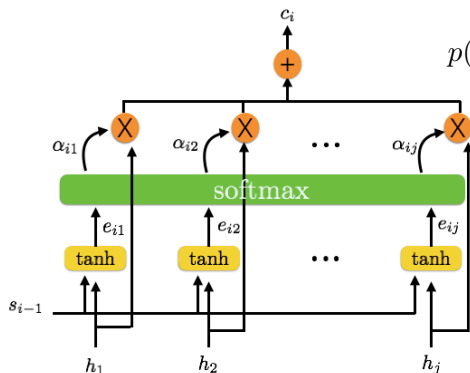
[Cited by 1276](#) [Related articles](#) [All 12 versions](#) [Cite](#) [Save](#)



Bidirectional RNN for Annotating Sequence



Bidirectional RNN for Annotating Sequence



$$p(y_o | y_1, \dots, y_{o-1}, c_o) = g(y_{o-1}, s_o, c_o)$$

$$s_o = f(y_{o-1}, s_{o-1}, c_o)$$

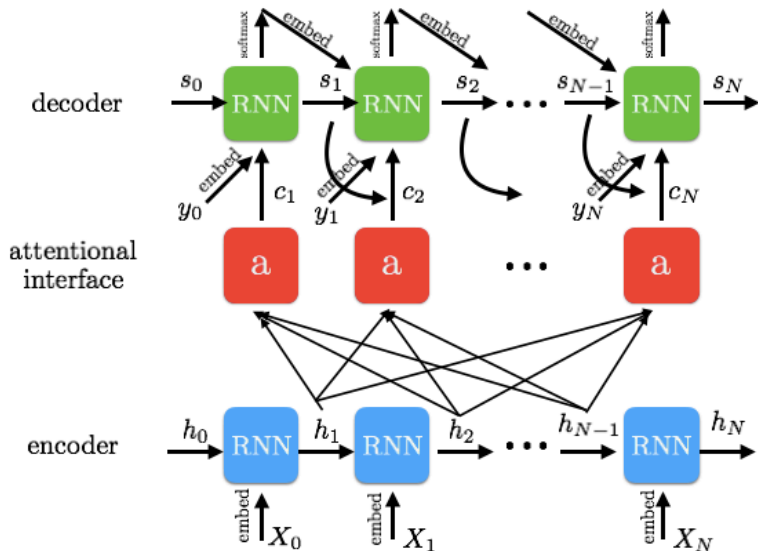
$$c_o = \sum_{t=1}^T \alpha_{ot} h_t$$

$$\alpha_{ot} = \frac{\exp(e_{ot})}{\sum_{k=1}^T \exp(e_{ok})}$$

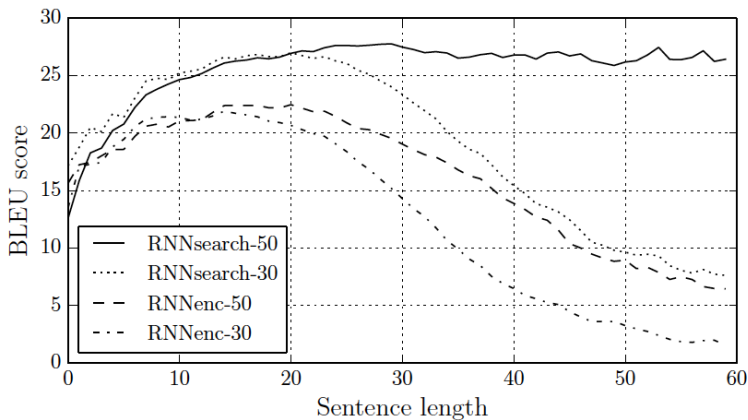
$$e_{ot} = r(s_{o-1}, h_t) \\ = v^T \tanh(W s_{o-1} + U h_t)$$

可见，计算顺序为： $s_{o-1} \rightarrow \alpha_{ot} \rightarrow c_o \rightarrow s_o$

Bidirectional RNN for Annotating Sequence



Bidirectional RNN for Annotating Sequence



Attention Roadmap - 2015

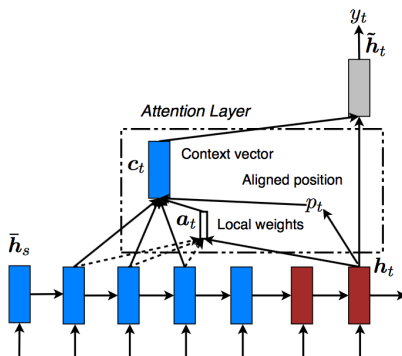
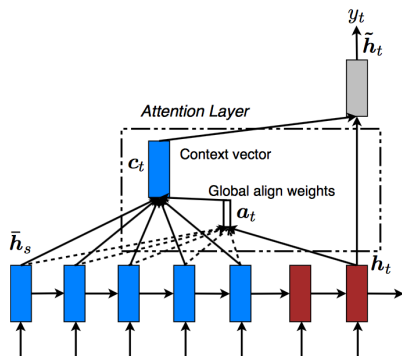
Global and Local Attentional Model

Effective approaches to attention-based neural machine translation

[MT Luong](#), [H Pham](#), [CD Manning](#) - arXiv preprint arXiv:1508.04025, 2015 - [arxiv.org](#)

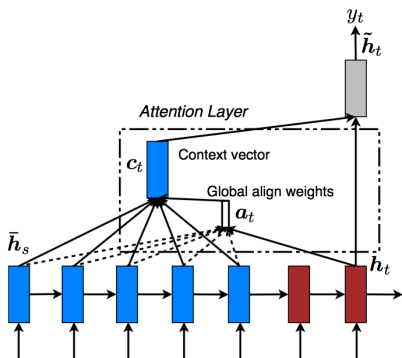
Abstract: An attentional mechanism has lately been used to improve neural machine translation (NMT) by selectively focusing on parts of the source sentence during translation. However, there has been little work exploring useful architectures for **attention-based** NMT.

Cited by 240 Related articles All 20 versions Cite Save



Global and Local Attentional Model

Global Attentional Model



$$\tilde{h}_o = \tanh(W_c[c_o; s_o])$$

$$p(y_o|y_1, \dots, y_{o-1}, x) = \text{softmax}(W_s \tilde{h}_o)$$

$$c_o = \sum_t h_t \alpha_{ot}$$

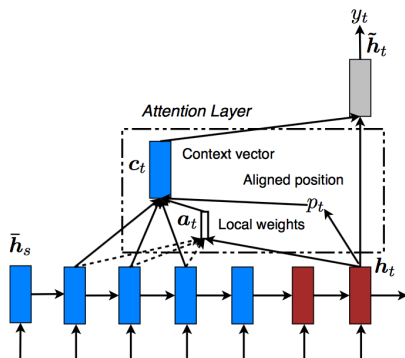
$$\begin{aligned} \alpha_{ot} &= \text{align}(s_o, h_t) \\ &= \frac{\exp(\text{score}(s_o, h_t))}{\sum_{t'}^T \exp(\text{score}(s_o, h_{t'}))} \end{aligned}$$

$$\text{score}(s_o, h_t) = \begin{cases} s_o^T h_t \\ s_o^T W_a h_t \\ v_a^T \tanh(W_a[s_o; h_t]) \end{cases}$$

可见，计算顺序为： $s_o \rightarrow \alpha_{ot} \rightarrow c_o \rightarrow \tilde{h}_o$

Global and Local Attentional Model

Local Attentional Model



在一个窗口中计算上下文信息 c_t ，但关键是如何选取窗口：

(1) 指定一个窗口: $[p_t - D, p_t + D]$ ，其中 $p_t = t$ ， D 为经验参数；

(2) 计算一个窗口：

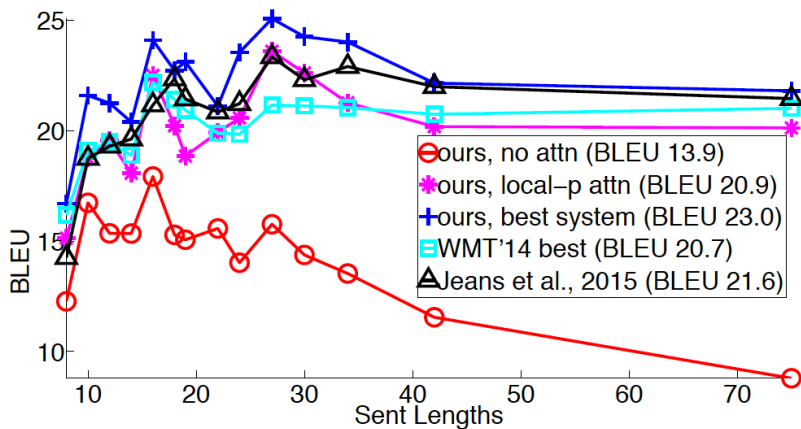
$$p_t = S \text{sigmoid}(v_p^T) \tanh(W_p h_t)$$

其中， W_p 与 v_p 为模型参数， S 为源句子长度；

(3) 更进一步，以 p_t 为中心，对 α_{ot} 做高斯：

$$\alpha_{ot} = \text{align}(s_o, h_t) \exp\left(-\frac{(s - p_t)^2}{2\delta^2}\right)$$

Global and Local Attentional Model



Thanks.