# Deep learning for NLP

Models of meaning

# Why Deep learning?

# Let's start with word vectors

1. The first try: one-hot encoding

```
V = {zebra, horse, school, summer}

      v(zebra)  = [1, 0, 0, 0]
      v(horse)  = [0, 1, 0, 0]
     v(school)  = [0, 0, 1, 0]
     v(summer)  = [0, 0, 0, 1]
```
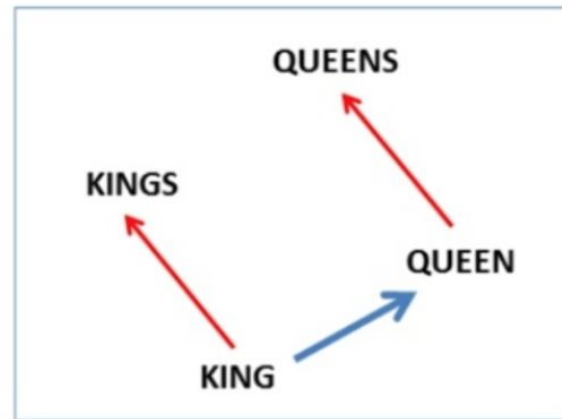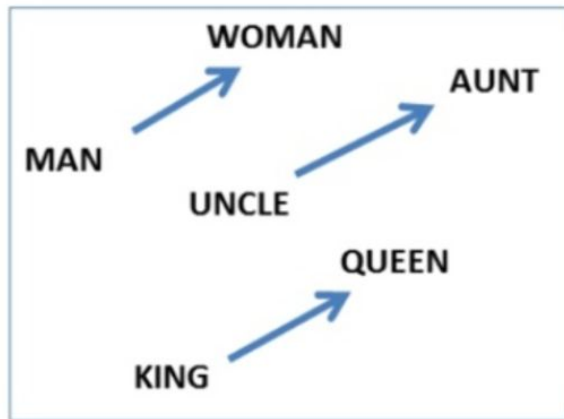
# The second try: word2vec and GloVe

$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\|\|\mathbf{B}\|} = \frac{\sum\limits_{i=1}^{n} A_i B_i}{\sqrt{\sum\limits_{i=1}^{n} A_i^2}\sqrt{\sum\limits_{i=1}^{n} B_i^2}}$$
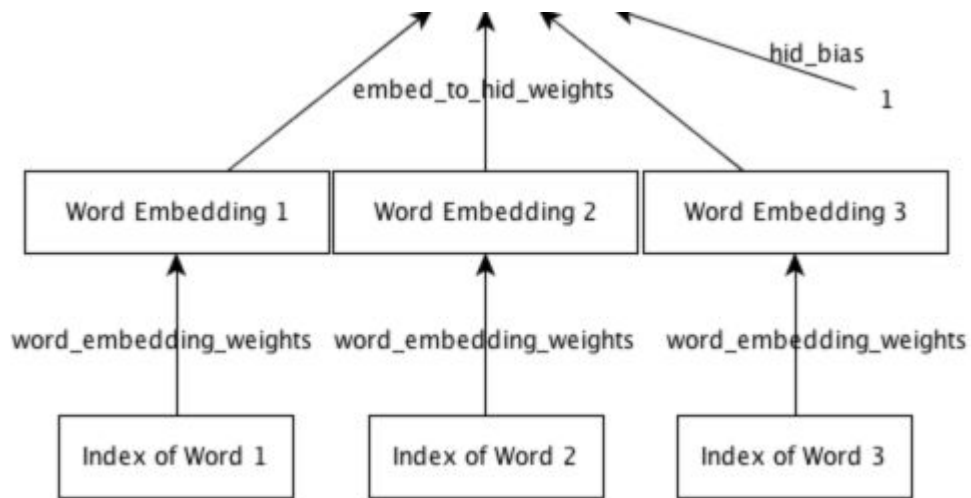
vec("man") – vec("king") + vec("woman") = vec("queen")

# Embedding layer

Embedding matrix:

$$L = \begin{bmatrix} \bullet & \bullet & \bullet & & \bullet & \bullet \\ \bullet & \bullet & \bullet & \cdots & \bullet & \bullet \\ \bullet & \bullet & \bullet & & \bullet & \bullet \\ \bullet & \bullet & \bullet & & \bullet & \bullet \end{bmatrix} n$$

|V|

the  cat   mat ...

# Feedforward Neural Networks for classification



$i$-th output $= P(w_t = i \,|\, context)$

softmax

most computation here

tanh

$C(w_{t-n+1})$     $C(w_{t-2})$   $C(w_{t-1})$

Table look-up in $C$

Matrix $C$ shared parameters across words

index for $w_{t-n+1}$     index for $w_{t-2}$     index for $w_{t-1}$

# Training Neural Networks

Indexes of words

↓

Embedding layer

↓

Neural Network

↓

Output (Softmax)

↓

Loss function (J)

↓

SGD on mini-batches

$$\theta^{new} = \theta^{old} - \alpha \nabla_\theta J(\theta)$$

# Convolution and pooling operations

Center element of the kernel is placed over the source pixel. The source pixel is then replaced with a weighted sum of itself and nearby pixels.

Source pixel

Convolution kernel
(emboss)

New pixel value (destination pixel)

(4 x 0)
(0 x 0)
(0 x 0)
(0 x 0)
(0 x 1)
(0 x 1)
(0 x 0)
(0 x 1)
+ (-4 x 2)
-8

## Single depth slice

x

| 1 | 1 | 2 | 4 |
|---|---|---|---|
| 5 | 6 | 7 | 8 |
| 3 | 2 | 1 | 0 |
| 1 | 2 | 3 | 4 |

y

max pool with 2x2 filters
and stride 2

| 6 | 8 |
|---|---|
| 3 | 4 |

# 1D/2D Convolutional neural networks

Text is not a picture! We have to convolve over words.



wait
for
the
video
and
do
n't
rent
it

| | | | |
| --- | --- | --- | --- |
| $n \times k$ representation of sentence with static and non-static channels | Convolutional layer with multiple filter widths and feature maps | Max-over-time pooling | Fully connected layer with dropout and softmax output |

# Dynamic Convolutional Neural Network



There are three tricky operation here:
- Convolution along rows for 2D filters
- Dynamic k-max pooling
- Folding: sum two rows in one

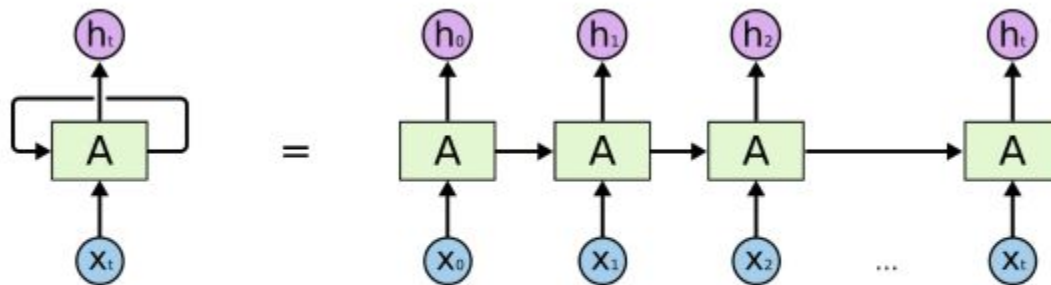k for every layer is calculated by the following formula:

$$k_l = \max\left( k_{top},\ \lceil \frac{L-l}{L}s \rceil \right)$$

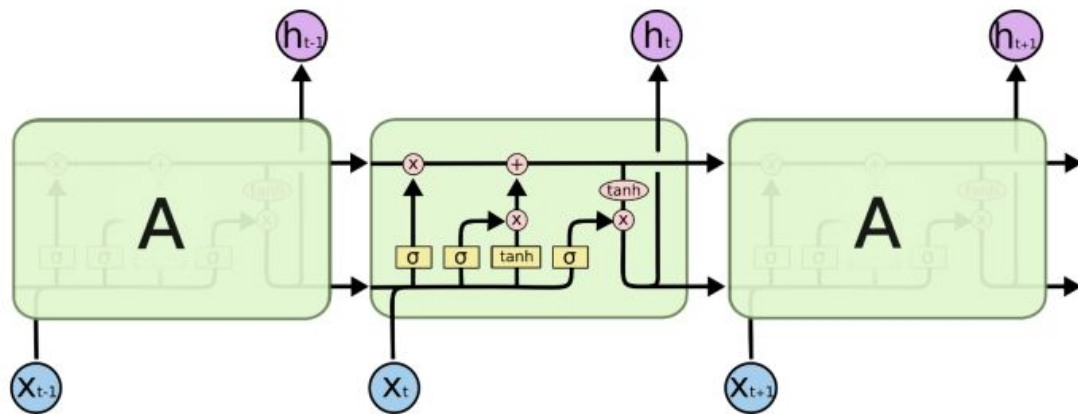http://arxiv.org/pdf/1404.2188v1.pdf

# RNN 1: Vanilla RNN

$$h_t = \tanh(W_h h_{t-1} + W_x x_t),$$



An unrolled recurrent neural network.

# RNN 2: LSTM



The repeating module in an LSTM contains four interacting layers.

$$f_t = \sigma\left(W_f \cdot [h_{t-1}, x_t] + b_f\right)$$
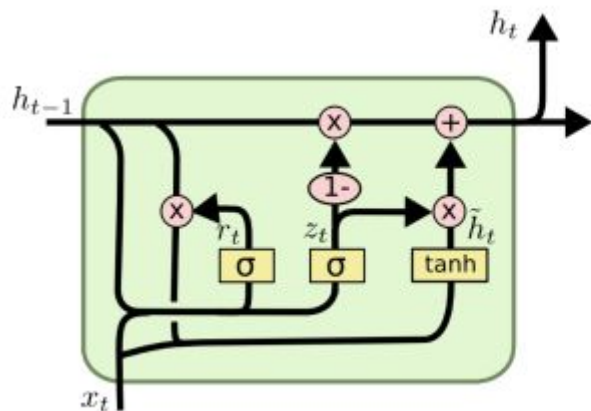
$$i_t = \sigma\left(W_i \cdot [h_{t-1}, x_t] + b_i\right)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

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

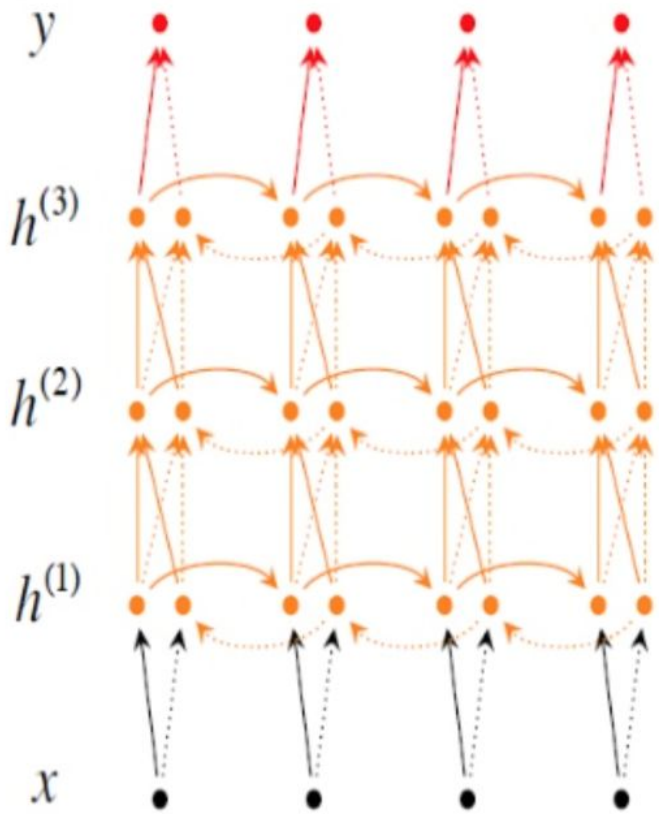$$o_t = \sigma\left(W_o\ [h_{t-1}, x_t] + b_o\right)$$

$$h_t = o_t * \tanh\left(C_t\right)$$

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

# RNN 3: GRU



$$z_t = \sigma \left( W_z \cdot [h_{t-1}, x_t] \right)$$

$$r_t = \sigma \left( W_r \cdot [h_{t-1}, x_t] \right)$$

$$\tilde{h}_t = \tanh \left( W \cdot [r_t * h_{t-1}, x_t] \right)$$

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

# Multi-layer bidirectional RNN



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

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

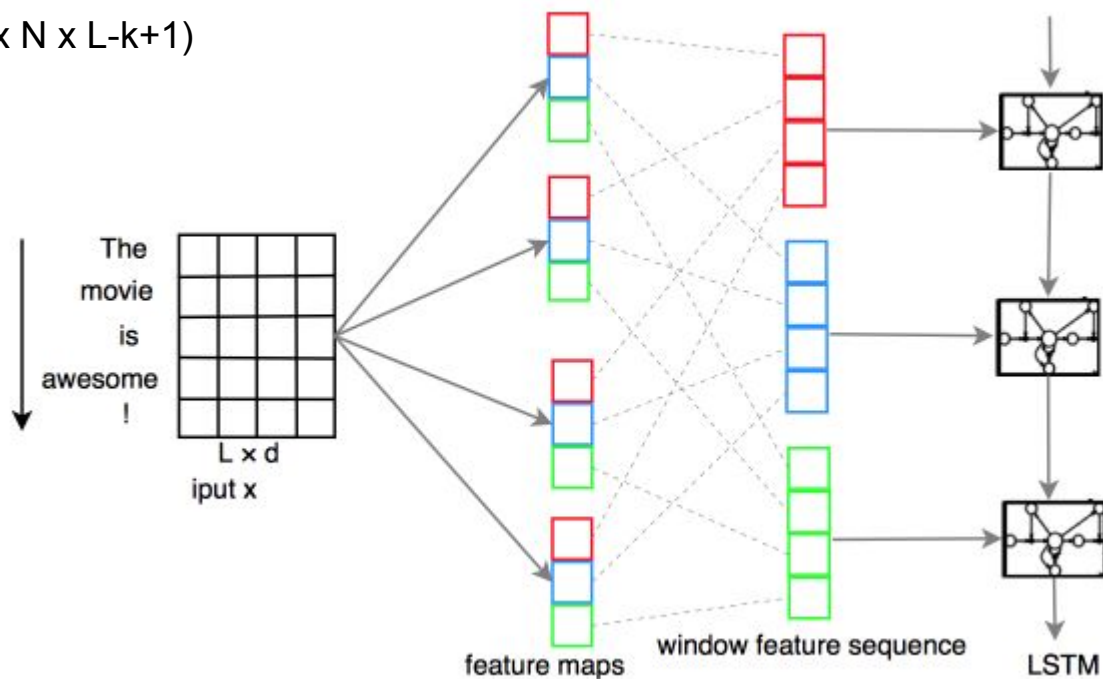$$y_t = g(U[\overrightarrow{h_t}^{(L)}; \overleftarrow{h_t}^{(L)}] + c)$$

# C-LSTM Neural Networks

After CNN we have a tensor: (B x N x L-k+1)

, where B - size of mini-batch
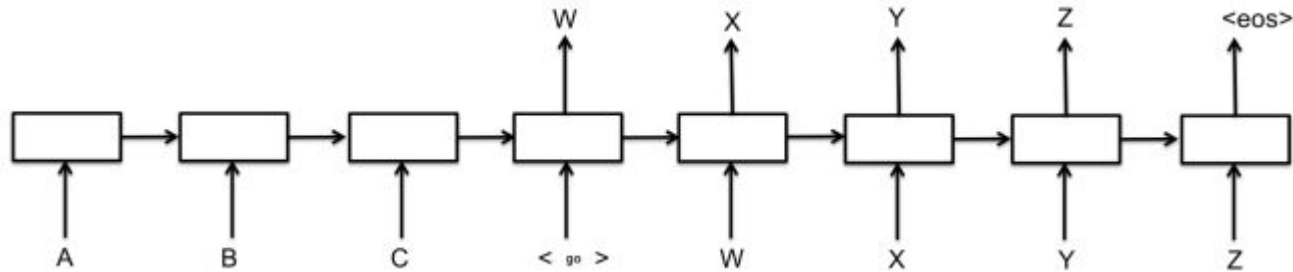N - count of filters
L - max length of sentence
k - length of filters



The
movie
is
awesome
!

L × d
iput x

feature maps        window feature sequence        LSTM

# Neural machine translation (NMT)



$$1/|\mathcal{S}| \sum_{(T,S)\in\mathcal{S}} \log p(T|S)$$

A,B,C – vectors of words of one language
W,X,Y,Z (below) – vectors of words of another language
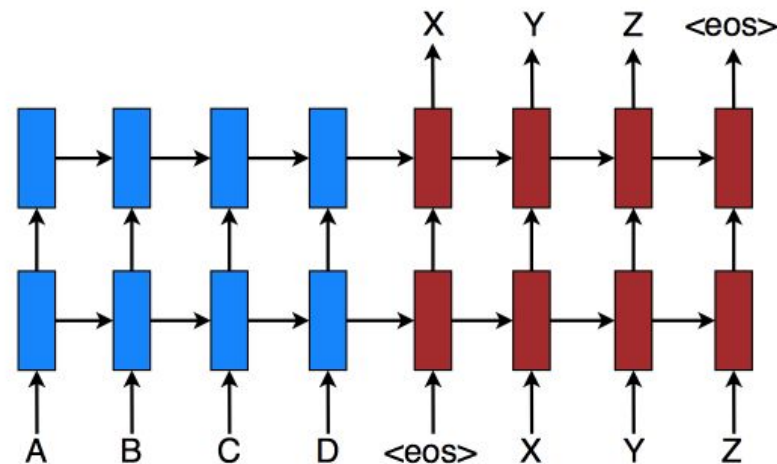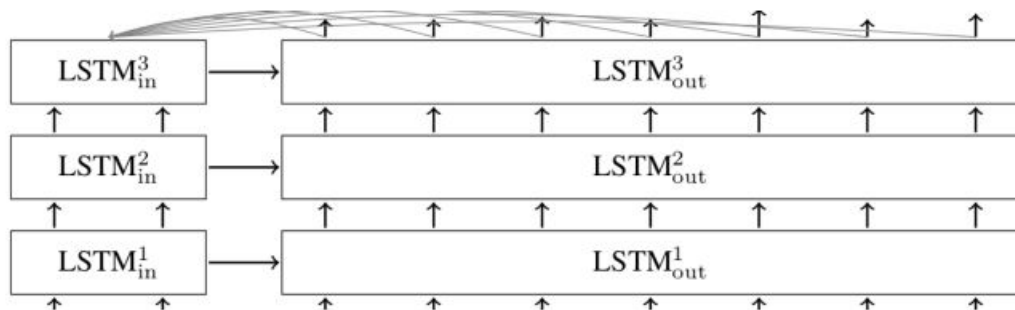W,X,Y,Z (above) – one-hot vectors for another language
<go>, <EOS> – special vectors of the start and end of output sentence
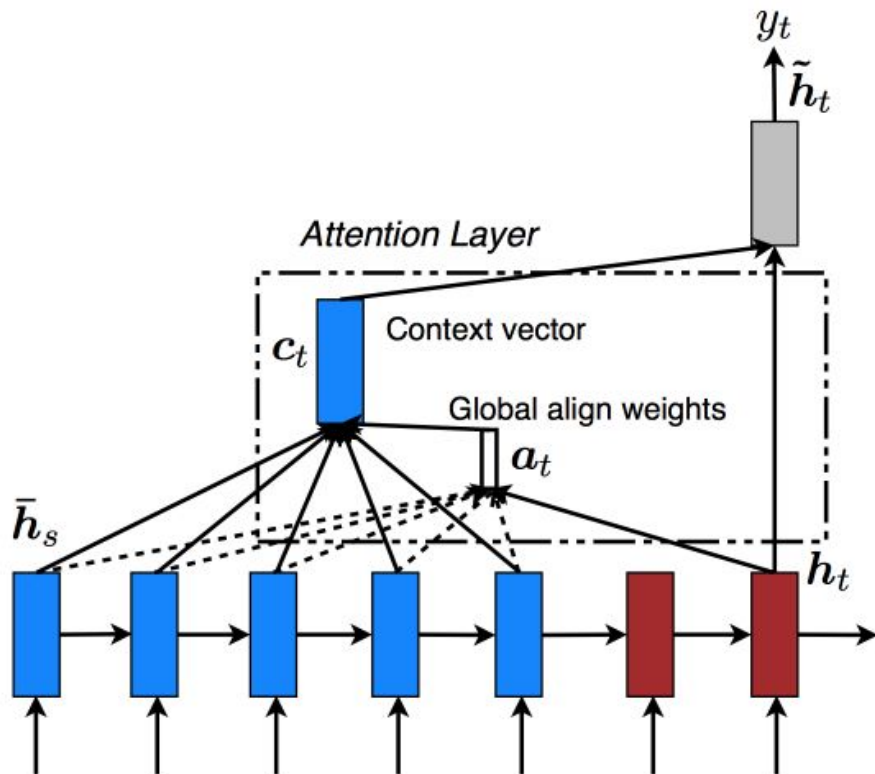Every rectangular block – LSTM block

# NMT with several layers



The order of the words of the input sentence was reversed!

# NMT + ATTENTION



$$\text{score}(\boldsymbol{h}_t, \bar{\boldsymbol{h}}_s) = \begin{cases} \boldsymbol{h}_t^\top \bar{\boldsymbol{h}}_s & dot \\ \boldsymbol{h}_t^\top \boldsymbol{W_a} \bar{\boldsymbol{h}}_s & general \\ \boldsymbol{W_a}[\boldsymbol{h}_t; \bar{\boldsymbol{h}}_s] & concat \end{cases}$$

$$\boldsymbol{a}_t(s) = \text{align}(\boldsymbol{h}_t, \bar{\boldsymbol{h}}_s)$$

$$= \frac{\exp\left(\text{score}(\boldsymbol{h}_t, \bar{\boldsymbol{h}}_s)\right)}{\sum_{s'} \exp\left(\text{score}(\boldsymbol{h}_t, \bar{\boldsymbol{h}}_{s'})\right)}$$

A global context vector $c_t$ is then computed as the weighted average, according to $a_t$, over all the source states.
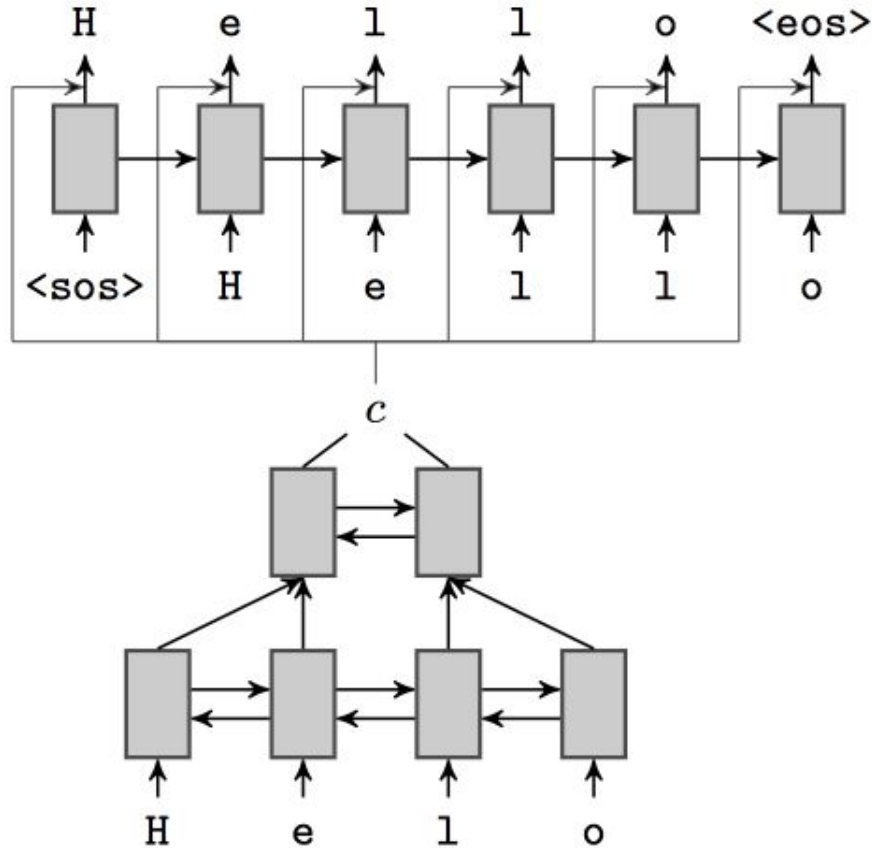
$$\tilde{\boldsymbol{h}}_t = \tanh(\boldsymbol{W_c}[\boldsymbol{c}_t; \boldsymbol{h}_t])$$

$$p(y_t|y_{<t}, x) = \text{softmax}(\boldsymbol{W_s} \tilde{\boldsymbol{h}}_t)$$

http://arxiv.org/pdf/1409.0473v6.pdf
http://www.aclweb.org/anthology/D15-1166

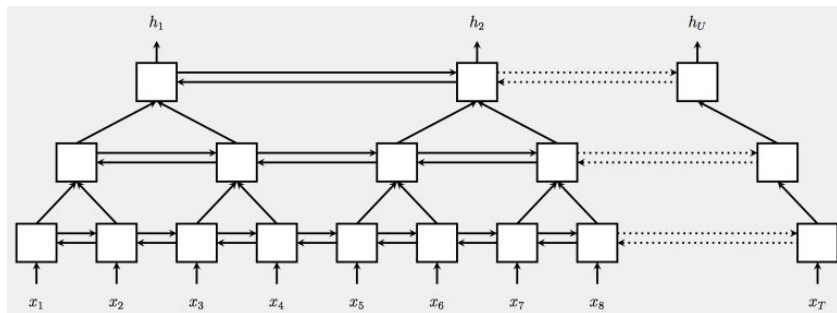# Character-level error correction with attention

# Formulas

## Encoder

$$f_t^{(j)} = \text{GRU}(f_{t-1}^{(j)}, c_t^{(j-1)}),$$

$$b_t^{(j)} = \text{GRU}(b_{t+1}^{(j)}, c_t^{(j-1)}),$$

$$h_t^{(j)} = f_t^{(j)} + b_t^{(j)}$$

$$c_t^{(j)} = \tanh\left(W_{\text{pyr}}^{(j)}\left[h_{2t}^{(j-1)}, h_{2t+1}^{(j-1)}\right]^\top + b_{\text{pyr}}^{(j)}\right)$$

## Decoder

$$d_t^{(j)} = \text{GRU}(d_{t-1}^{(j)}, d_t^{(j-1)}),$$

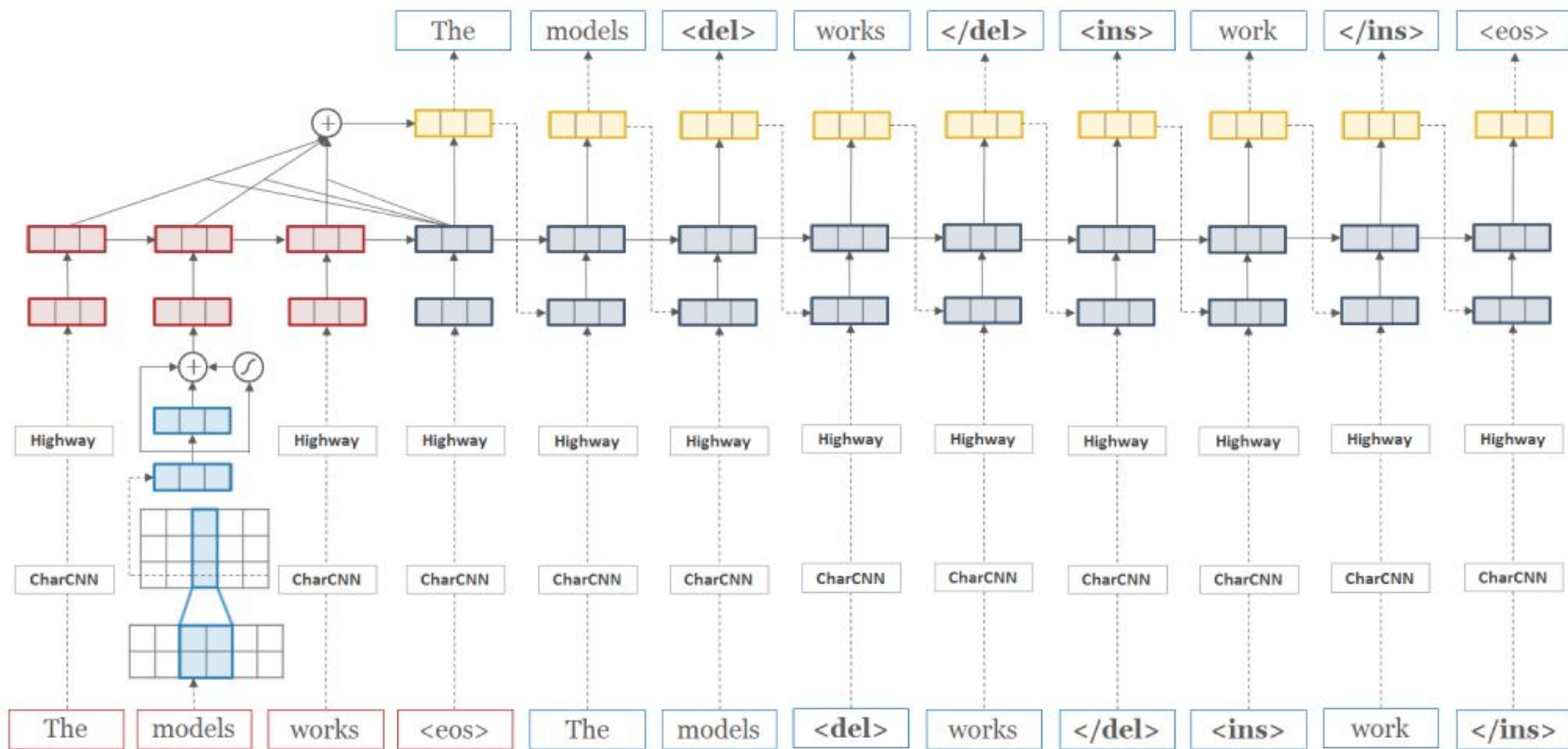### Attention

$$u_{tk} = \phi_1(d_t^{(M)})^\top \phi_2(c_k)$$

$$\alpha_{tk} = \frac{u_{tk}}{\sum_j u_{tj}}$$

$$a_t = \sum_j \alpha_{tj} c_j$$

### Output:

The weighted sum of the encoded hidden states **a**t is then concatenated with **d**(M), and passed through another affine transform followed by a ReLU nonlinearity before the final softmax output layer.

# Sentence-Level Grammatical Error Identification as Sequence-to-Sequence Correction



http://arxiv.org/pdf/1604.04677.pdf

# Formulas for word/character level

Encoder with attention to get context vector Cj:

Decoder

$$u_{j,i} = \mathbf{h}_j^t \cdot \mathbf{W}_\alpha \mathbf{h}_i^s$$

$$\alpha_{j,i} = \frac{\exp u_{j,i}}{\sum_{k \in [1,I]} \exp u_{j,k}}$$

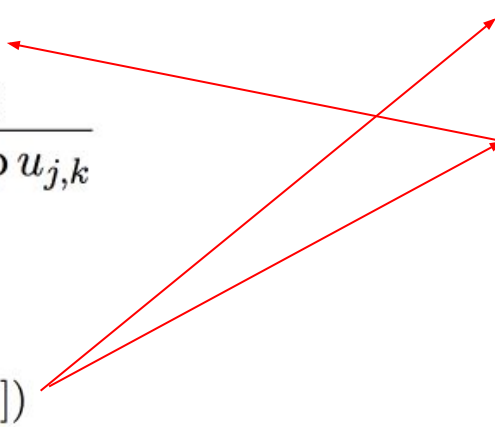$$\mathbf{v}_j = \sum_{i \in [1,I]} \alpha_{j,i} \mathbf{h}_i^s$$

$$\mathbf{c}_j = \tanh(\mathbf{W}[\mathbf{v}_j; \mathbf{h}_j^t])$$

$$p(t_{j+1} \mid \mathbf{s}, \mathbf{t}_{<j}) = \mathrm{softmax}(\mathbf{U}\mathbf{c}_j + \mathbf{b})$$

$$\mathbf{h}_j^t = \mathrm{LSTM}(\mathbf{h}_{j-1}^t, [\mathbf{x}_j^t; \mathbf{c}_{j-1}])$$

$$\mathbf{h}_0^t \leftarrow \mathbf{h}_I^s$$

# CharCNN and Highway Networks

Two separate CharCNNs for Encoder and Decoder:

Highway network:

$$\mathbf{f}_i[k] = \tanh(\langle \mathbf{P}_i[*, k : k + w - 1], \mathbf{H} \rangle + b)$$

$$\hat{\mathbf{z}} = \mathbf{r} \odot f(\mathbf{W}\mathbf{z} + \mathbf{b}) + (\mathbf{1} - \mathbf{r}) \odot \mathbf{z}$$

$$z_i = \max_k \mathbf{f}_i[k]$$

,where f is ReLu; r = σ(Wrz+br)

We use multiple filters H1,...Hh to obtain a vector zi ∈ Rh as the representation for a given source/target word or tag.