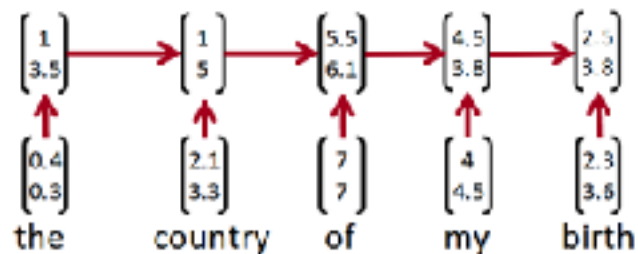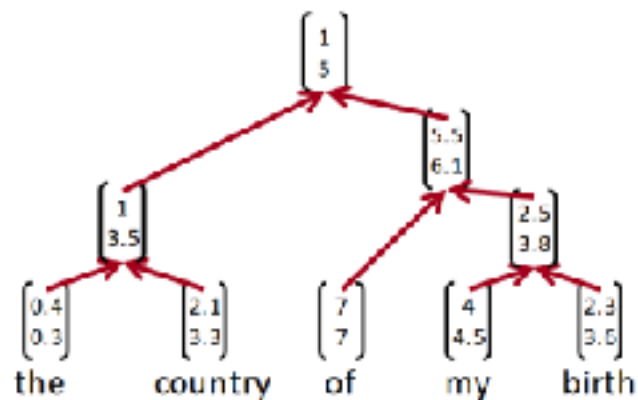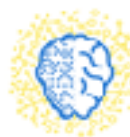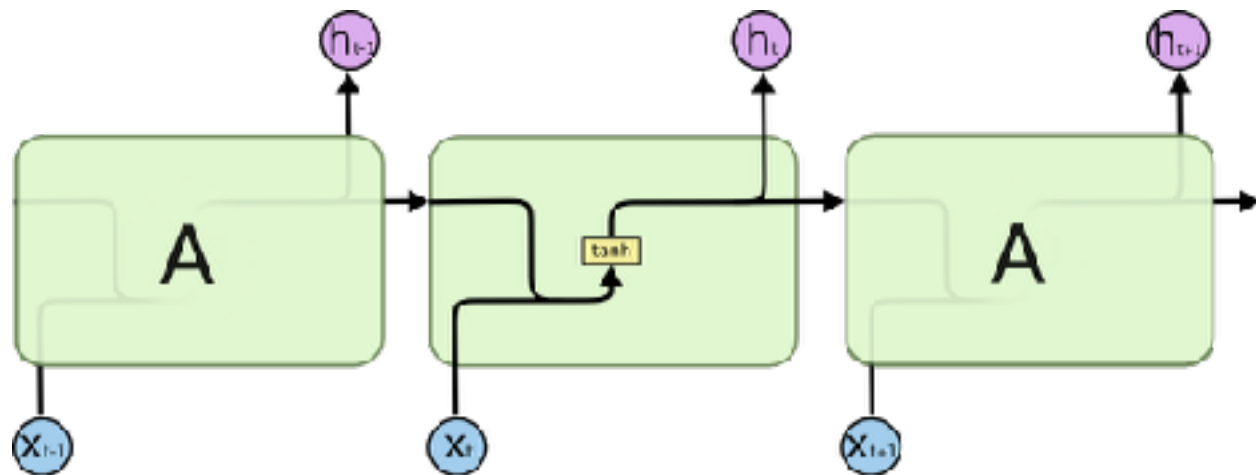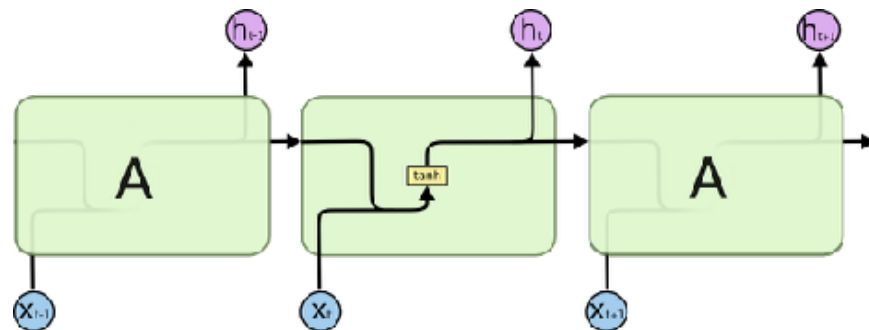# Recurrent Neural Networks part 2

Mikhail Arkhipov

Laboratory of Neural Systems and Deep Learning
MIPT

# Recursive vs. Recurrent

Stanford CS234d

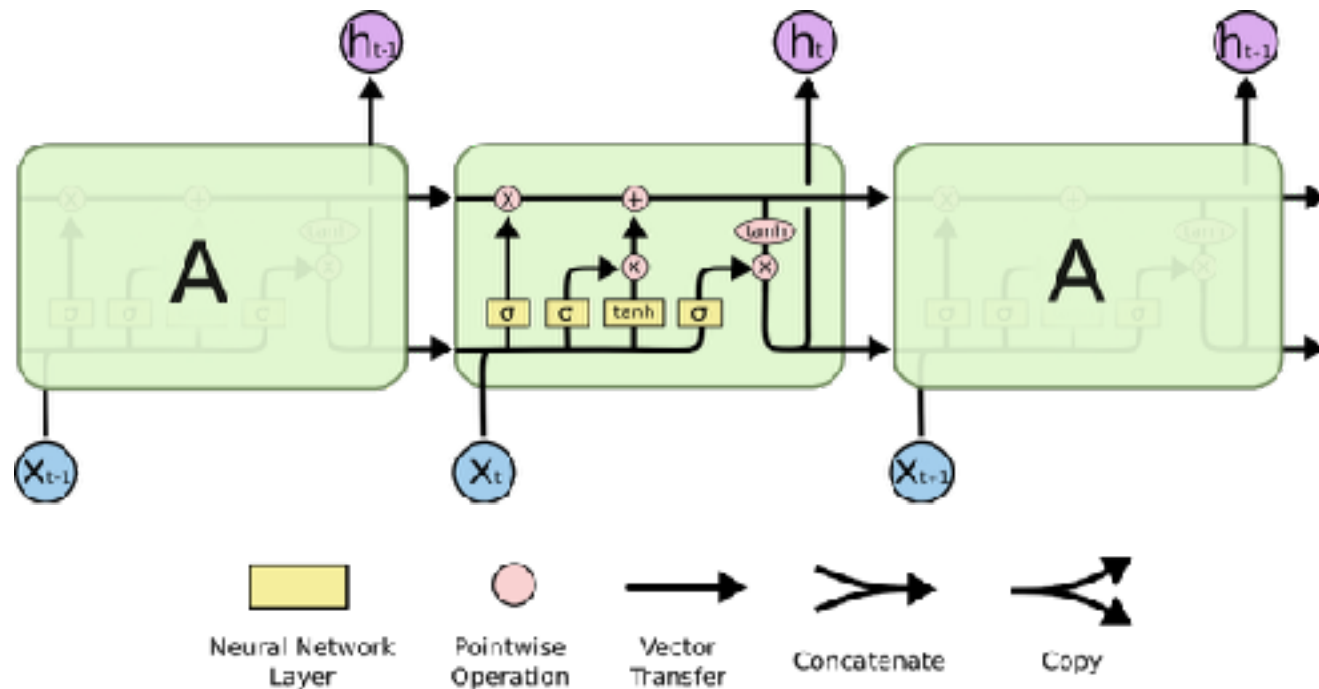- ● Exploding gradients
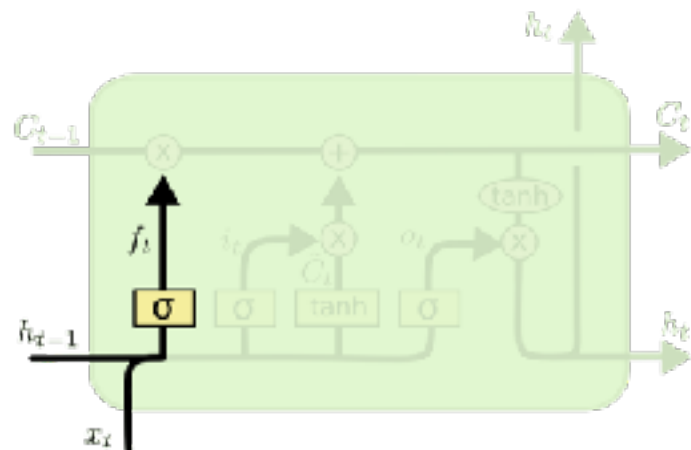- ● Vanishing gradients

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

# LSTM



Neural Network Layer · Pointwise Operation · Vector Transfer · Concatenate · Copy

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

[Hochreiter, Schmidhuber 1997]
Long Short-Term Memory

# LSTM

# LSTM



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

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

$$f_t = \sigma \left( W_f \cdot [h_{t-1}, x_t] \; + \; b_f \right)$$
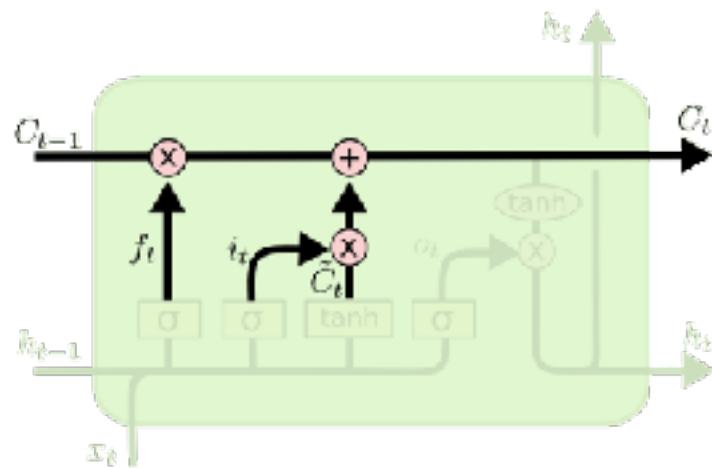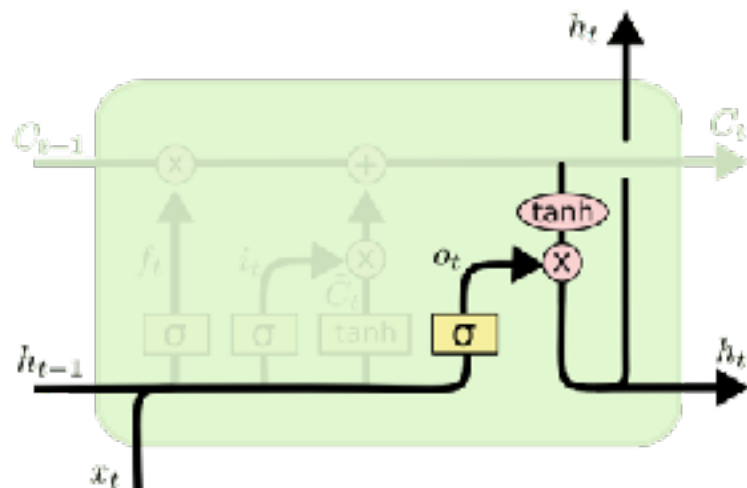
# LSTM



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

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

# LSTM



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

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

# LSTM



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

$$h_t = o_t * \tanh \left( C_t \right)$$
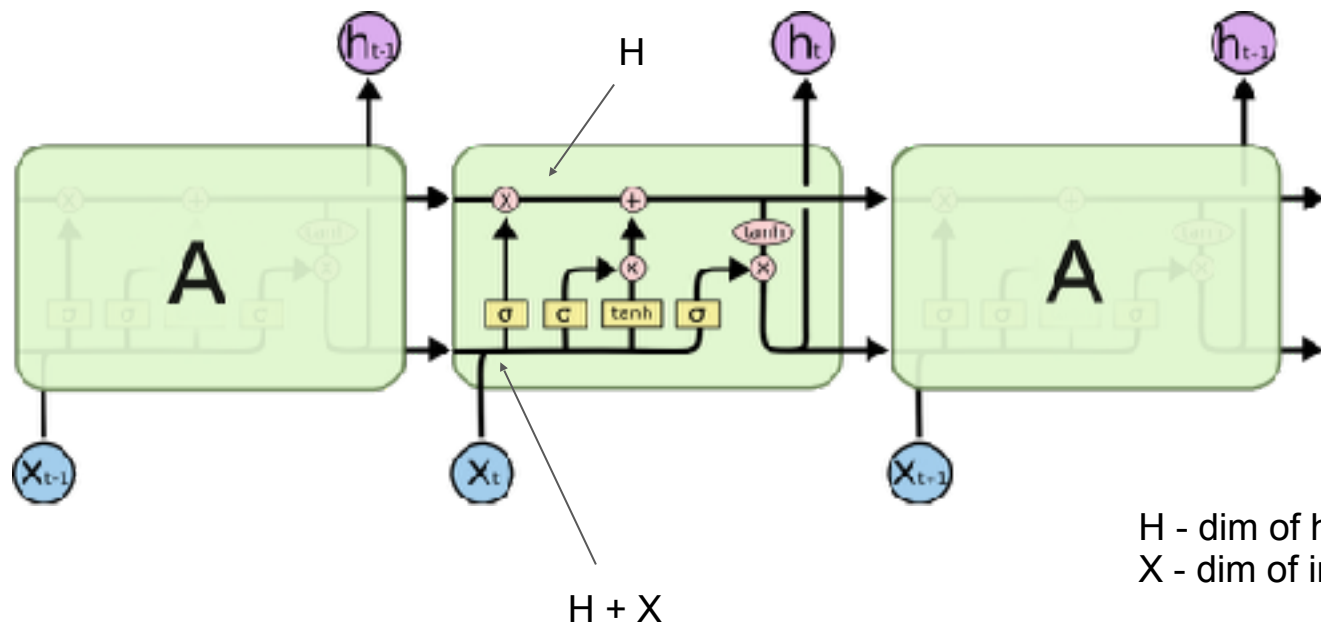
http://colah.github.io/posts/2015-08-Understanding-LSTMs/

# LSTM



H + X

H - dim of hidden state
X - dim of input features

# LSTM



H - dim of hidden state
X - dim of input features

H + X

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

# LSTM



dim of the cell state

H

H + X    dim of the input
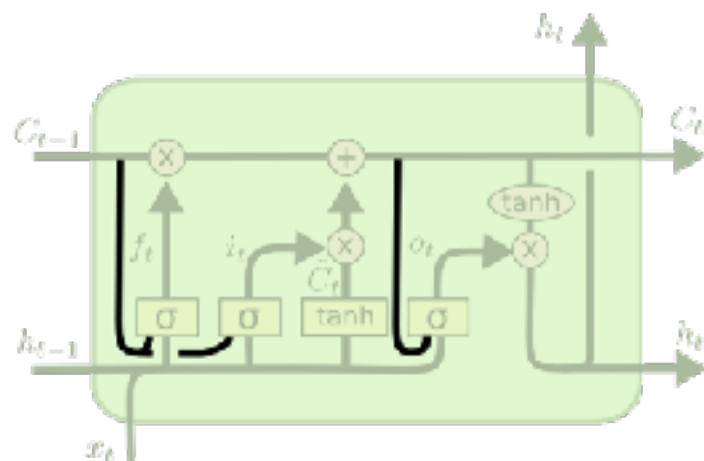
H - dim of hidden state
X - dim of input features

Total number of parameters   **(H + X) * H * 4**

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

# LSTM Coupled Gates



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

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

# LSTM Peephole connections



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

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

# LSTM overview



[Greff et. al 2017]
LSTM: A Search Space Odyssey
https://arxiv.org/pdf/1503.04069.pdf

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

# Gated Recurrent Unit (GRU)



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

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

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

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

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

(a) Standard Neural Net

(b) After applying dropout.

# Dropout



(a) Naive dropout RNN      (b) Variational RNN

State-of-the-Art Large Scale Language Modeling in 12 Hours With a Single GPU

# LayerNorm



$$\mathbf{h}^t = f\left[\frac{\mathbf{g}}{\sigma^t} \odot \left(\mathbf{a}^t - \mu^t\right) + \mathbf{b}\right] \qquad \mu^t = \frac{1}{H}\sum_{i=1}^{H} a_i^t \qquad \sigma^t = \sqrt{\frac{1}{H}\sum_{i=1}^{H}\left(a_i^t - \mu^t\right)^2}$$

# Autoregressive Recurrent models



$$p(\mathbf{x}) = \prod_{i} p(x|x_{<i}) = p(x_0)p(x_1|x_0)p(x_2|x_0, x_1)\ldots$$

$$p(\mathbf{x}) = \prod_i p(x|x_{<i}) = p(x_0)p(x_1|x_0)p(x_2|x_0, x_1)\ldots$$

# Language Models



$$p(l_1, l_2, \ldots, l_N) = \prod_{k=1}^{N} p(l_k \mid l_1, l_2, \ldots, l_{k-1}).$$

# Bi-directional Language Model

# ELMo



ELMo

$$R_k = \{\mathbf{x}_k^{LM}, \overrightarrow{\mathbf{h}}_{k,j}^{LM}, \overleftarrow{\mathbf{h}}_{k,j}^{LM} \mid j = 1, \ldots, L\}$$
$$= \{\mathbf{h}_{k,j}^{LM} \mid j = 0, \ldots, L\},$$

$$\mathbf{ELMo}_k^{task} = E(R_k; \Theta^{task}) = \gamma^{task} \sum_{j=0}^{L} s_j^{task} \mathbf{h}_{k,j}^{LM}.$$

iPavlov

# ELMo results

| TASK | PREVIOUS SOTA | | OUR BASELINE | ELMo + BASELINE | INCREASE (ABSOLUTE/ RELATIVE) |
|---|---|---|---|---|---|
| SQuAD | Liu et al. (2017) | 84.4 | 81.1 | 85.8 | 4.7 / 24.9% |
| SNLI | Chen et al. (2017) | 88.6 | 88.0 | $88.7 \pm 0.17$ | 0.7 / 5.8% |
| SRL | He et al. (2017) | 81.7 | 81.4 | 84.6 | 3.2 / 17.2% |
| Coref | Lee et al. (2017) | 67.2 | 67.2 | 70.4 | 3.2 / 9.8% |
| NER | Peters et al. (2017) | $91.93 \pm 0.19$ | 90.15 | $92.22 \pm 0.10$ | 2.06 / 21% |
| SST-5 | McCann et al. (2017) | 53.7 | 51.4 | $54.7 \pm 0.5$ | 3.3 / 6.8% |

# ELMo layer distribution

# seq2seq

TF seq2seq tutorial

# seq2seq

TF seq2seq tutorial

# Fusion



[Sriram et. al 17]

iPavlov.ai

Training

Inference

**Exposure bias!**

# Generating texts with RNN LM



Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

iPavlov.ai

# Generating texts with RNN LM



A. Karpathy. The Unreasonable Effectiveness of Recurrent Neural Networks

Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

# Searching for interpretable cells



Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

# Unsupervised sentiment neuron

# Training algorithms

Spasibo