

Recurrent Neural Networks II

Milan Straka

 April 15, 2019



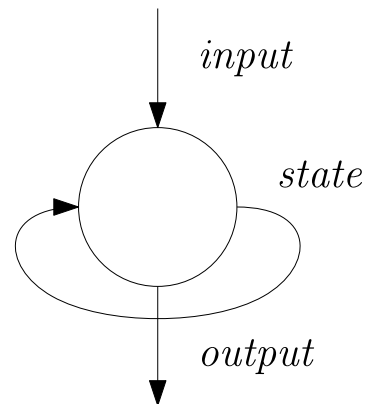
EUROPEAN UNION
European Structural and Investment Fund
Operational Programme Research,
Development and Education

Charles University in Prague
Faculty of Mathematics and Physics
Institute of Formal and Applied Linguistics

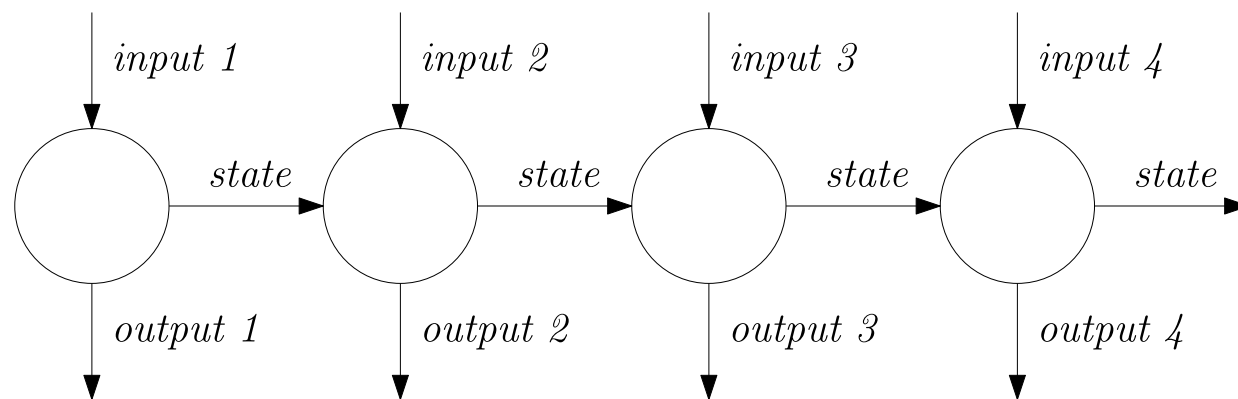


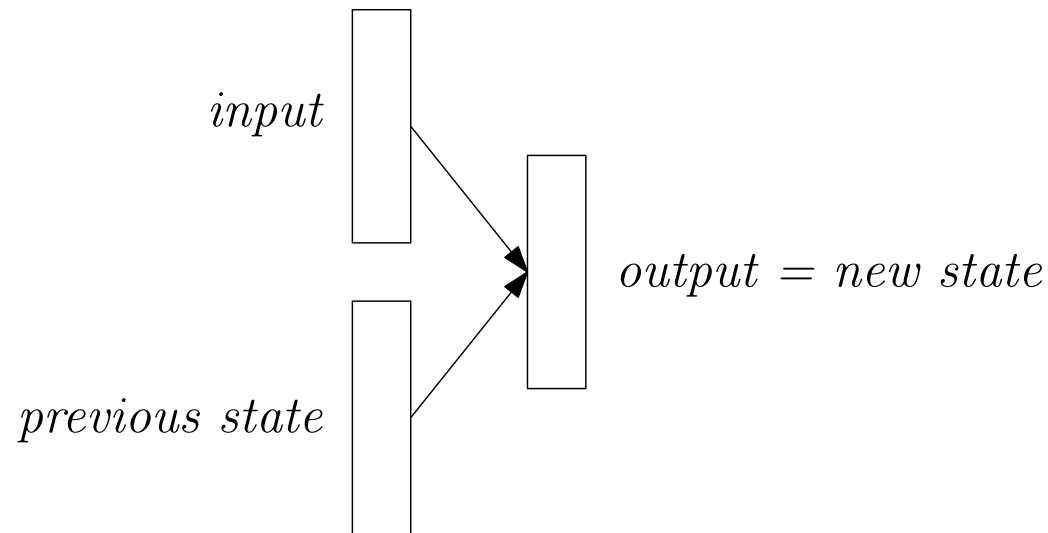
unless otherwise stated

Single RNN cell



Unrolled RNN cells





Given an input $\mathbf{x}^{(t)}$ and previous state $\mathbf{s}^{(t-1)}$, the new state is computed as

$$\mathbf{s}^{(t)} = f(\mathbf{s}^{(t-1)}, \mathbf{x}^{(t)}; \boldsymbol{\theta}).$$

One of the simplest possibilities is

$$\mathbf{s}^{(t)} = \tanh(\mathbf{U}\mathbf{s}^{(t-1)} + \mathbf{V}\mathbf{x}^{(t)} + \mathbf{b}).$$

Basic RNN cells suffer a lot from vanishing/exploding gradients (*the challenge of long-term dependencies*).

If we simplify the recurrence of states to

$$\mathbf{s}^{(t)} = \mathbf{U} \mathbf{s}^{(t-1)},$$

we get

$$\mathbf{s}^{(t)} = \mathbf{U}^t \mathbf{s}^{(0)}.$$

If \mathbf{U} has eigenvalue decomposition of $\mathbf{U} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^{-1}$, we get

$$\mathbf{s}^{(t)} = \mathbf{Q} \mathbf{\Lambda}^t \mathbf{Q}^{-1} \mathbf{s}^{(0)}.$$

The main problem is that the *same* function is iteratively applied many times.

Several more complex RNN cell variants have been proposed, which alleviate this issue to some degree, namely **LSTM** and **GRU**.

Long Short-Term Memory

Later in Gers, Schmidhuber & Cummins (1999) a possibility to *forget* information from memory cell c_t was added.

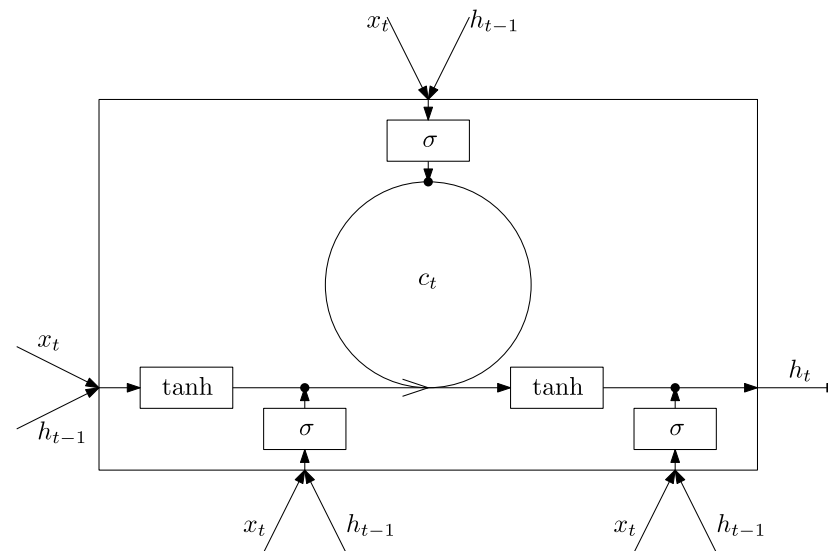
$$i_t \leftarrow \sigma(W^i x_t + V^i h_{t-1} + b^i)$$

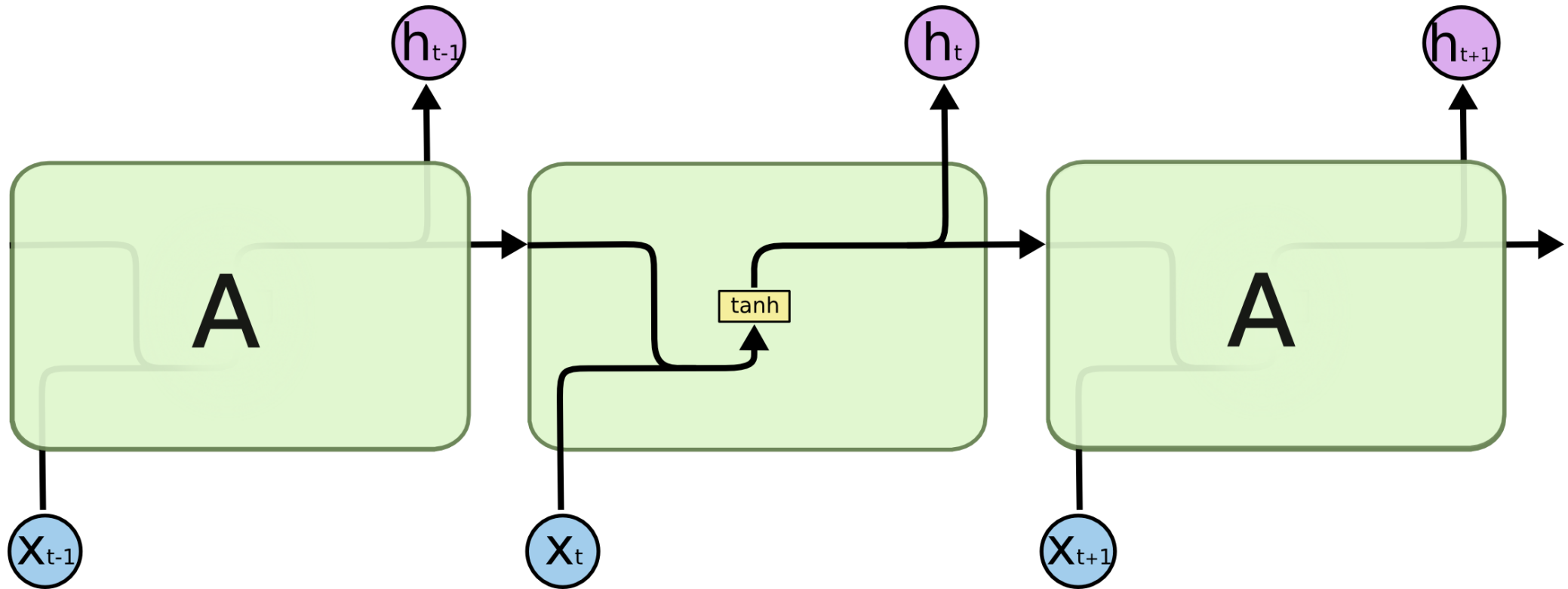
$$f_t \leftarrow \sigma(W^f x_t + V^f h_{t-1} + b^f)$$

$$o_t \leftarrow \sigma(W^o x_t + V^o h_{t-1} + b^o)$$

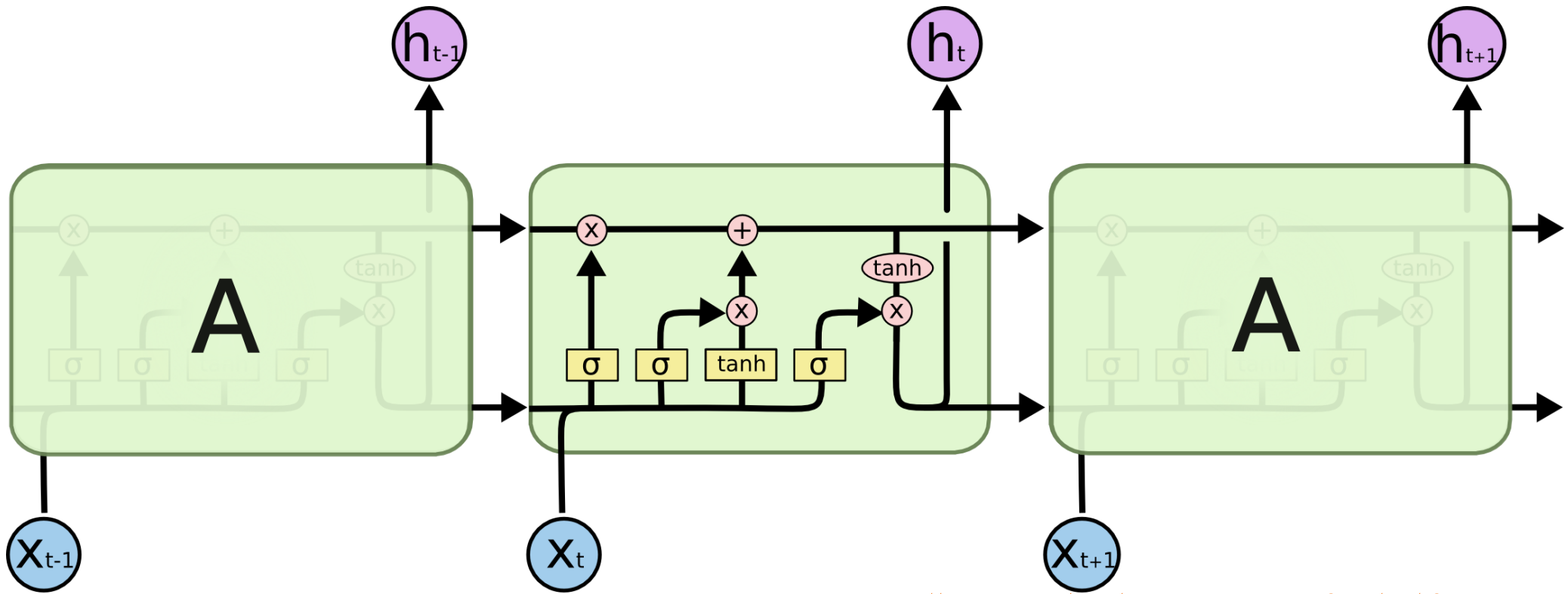
$$c_t \leftarrow f_t \cdot c_{t-1} + i_t \cdot \tanh(W^y x_t + V^y h_{t-1} + b^y)$$

$$h_t \leftarrow o_t \cdot \tanh(c_t)$$



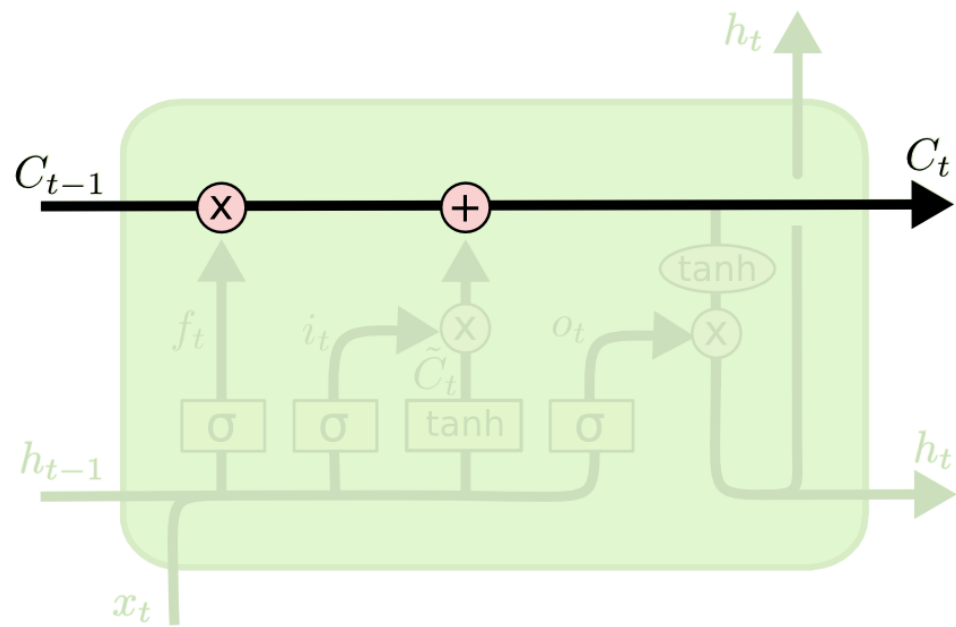


<http://colah.github.io/posts/2015-08-Understanding-LSTMs/img/LSTM3-SimpleRNN.png>



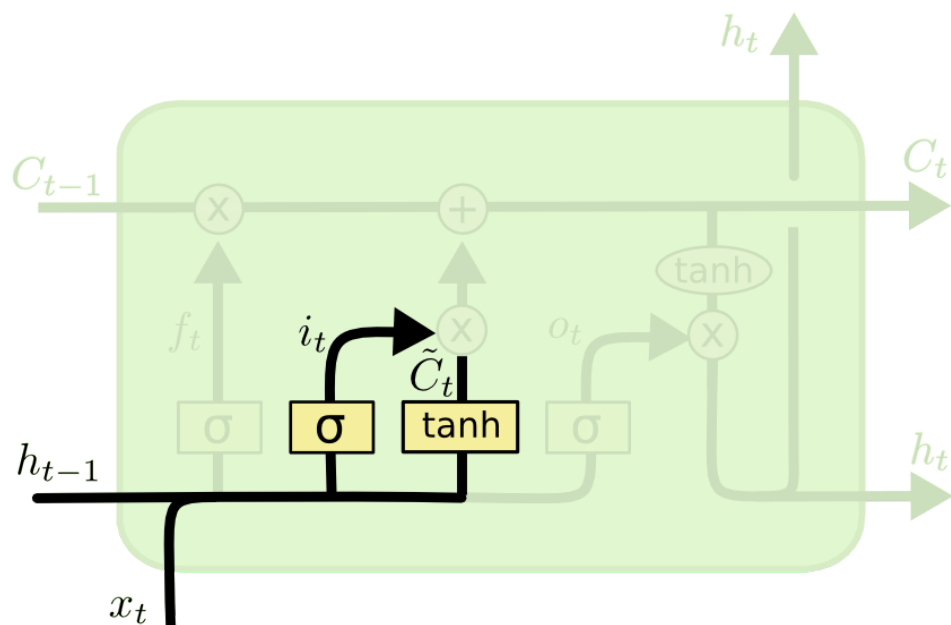
<http://colah.github.io/posts/2015-08-Understanding-LSTMs/img/LSTM3-chain.png>

Long Short-Term Memory



<http://colah.github.io/posts/2015-08-Understanding-LSTMs/img/LSTM3-C-line.png>

Long Short-Term Memory

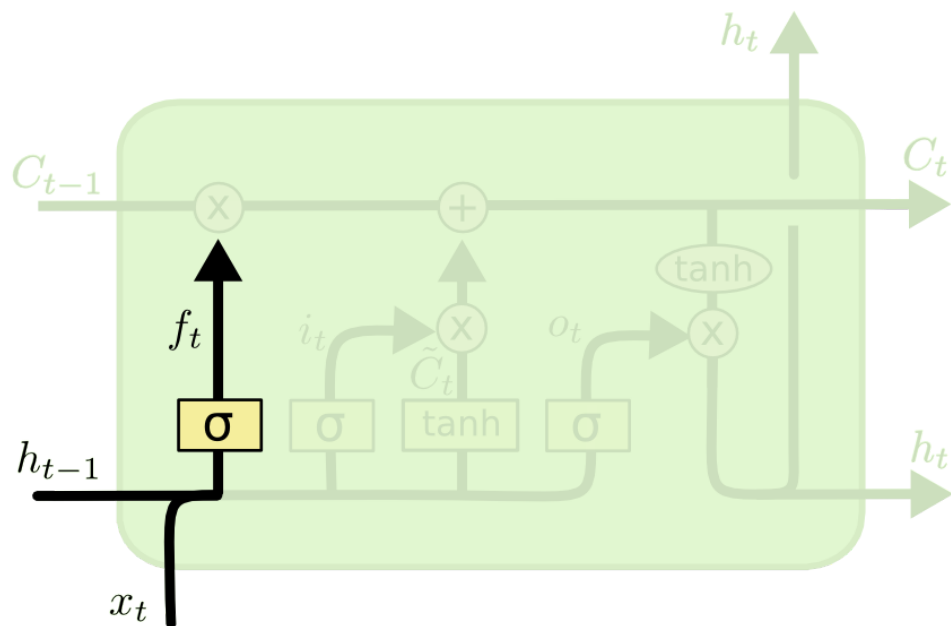


$$i_t = \sigma (W_i \cdot [h_{t-1}, x_t] + b_i)$$

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

<http://colah.github.io/posts/2015-08-Understanding-LSTMs/img/LSTM3-focus-i.png>

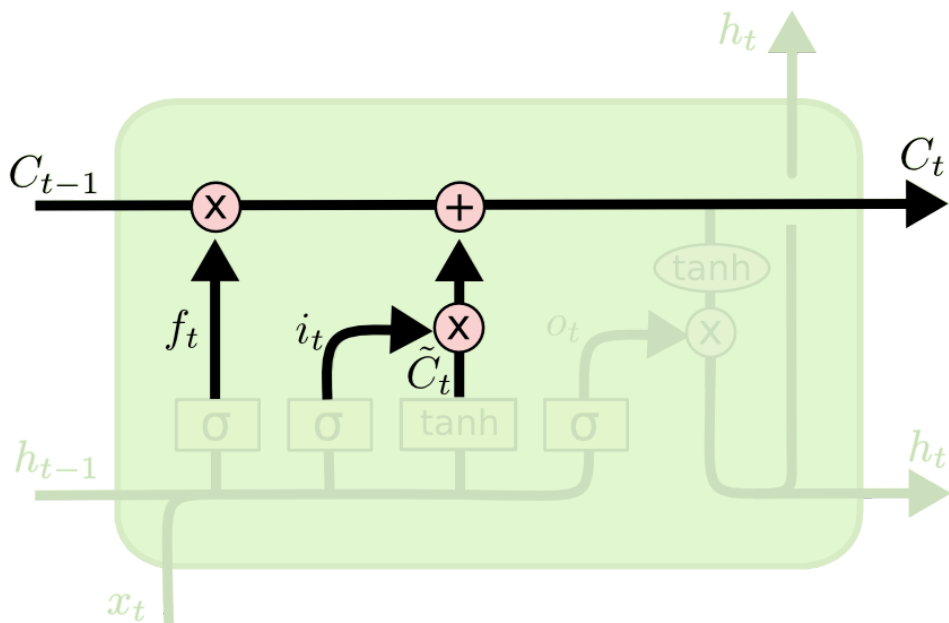
Long Short-Term Memory



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

<http://colah.github.io/posts/2015-08-Understanding-LSTMs/img/LSTM3-focus-f.png>

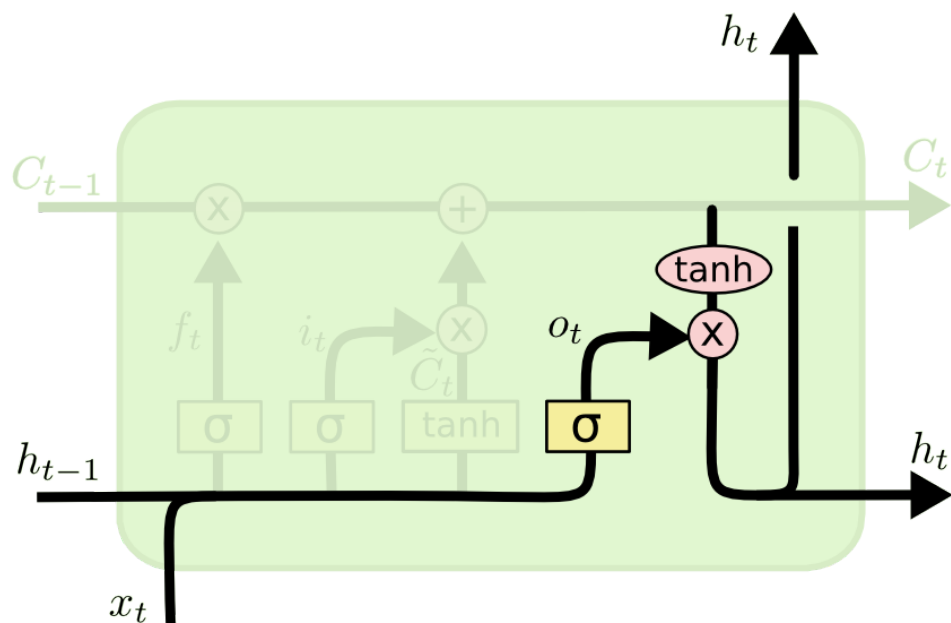
Long Short-Term Memory



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

<http://colah.github.io/posts/2015-08-Understanding-LSTMs/img/LSTM3-focus-C.png>

Long Short-Term Memory



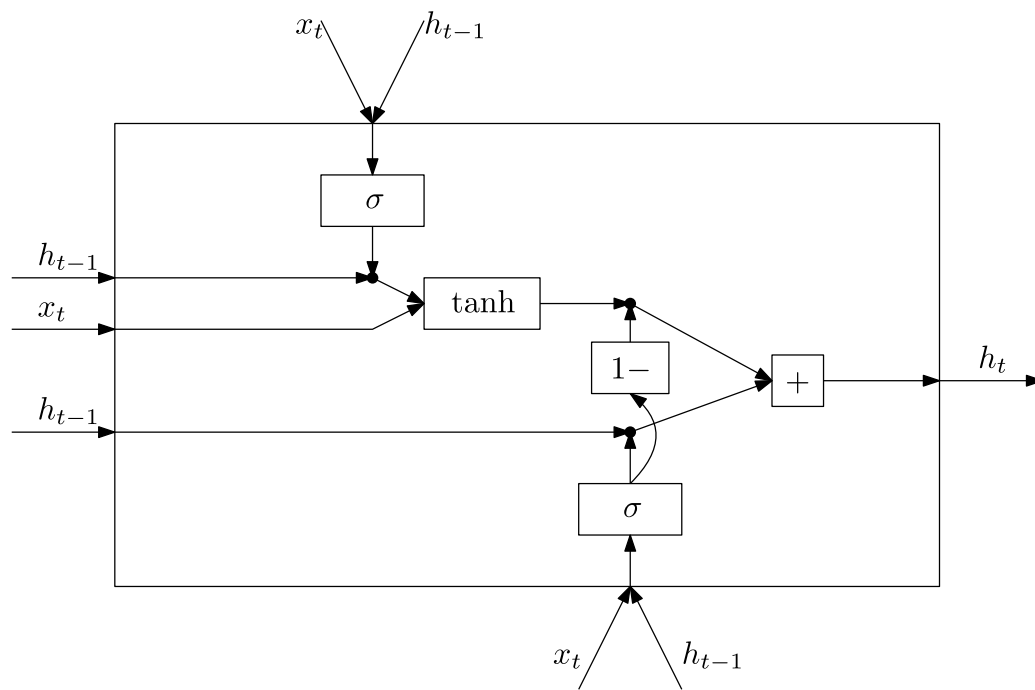
$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

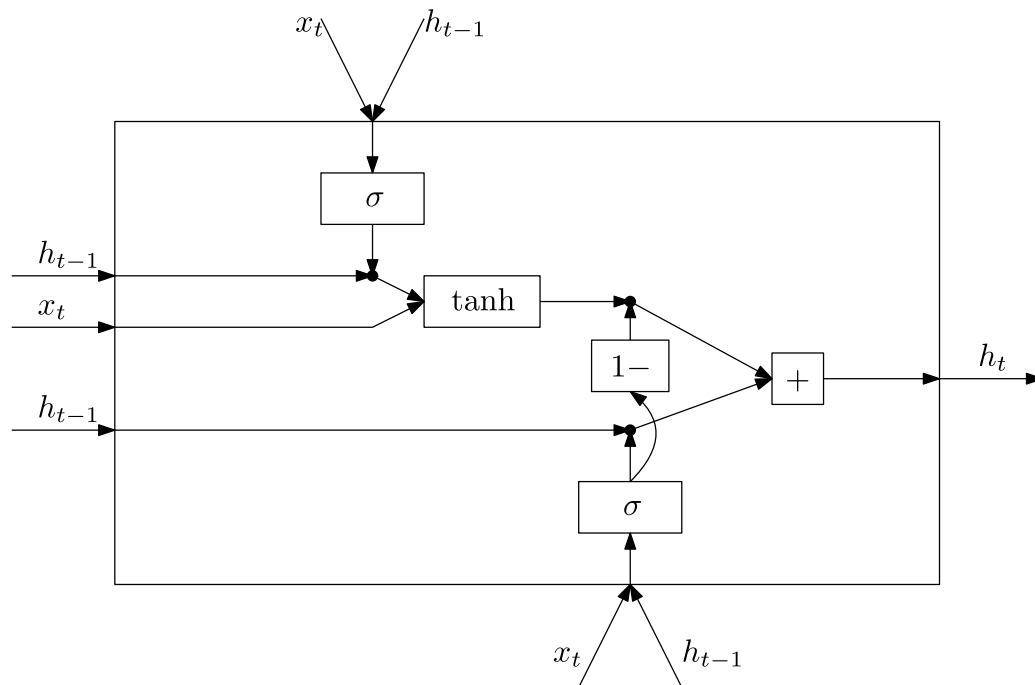
$$h_t = o_t * \tanh (C_t)$$

<http://colah.github.io/posts/2015-08-Understanding-LSTMs/img/LSTM3-focus-o.png>

Gated recurrent unit (GRU) was proposed by Cho et al. (2014) as a simplification of LSTM. The main differences are

- no memory cell
- forgetting and updating tied together



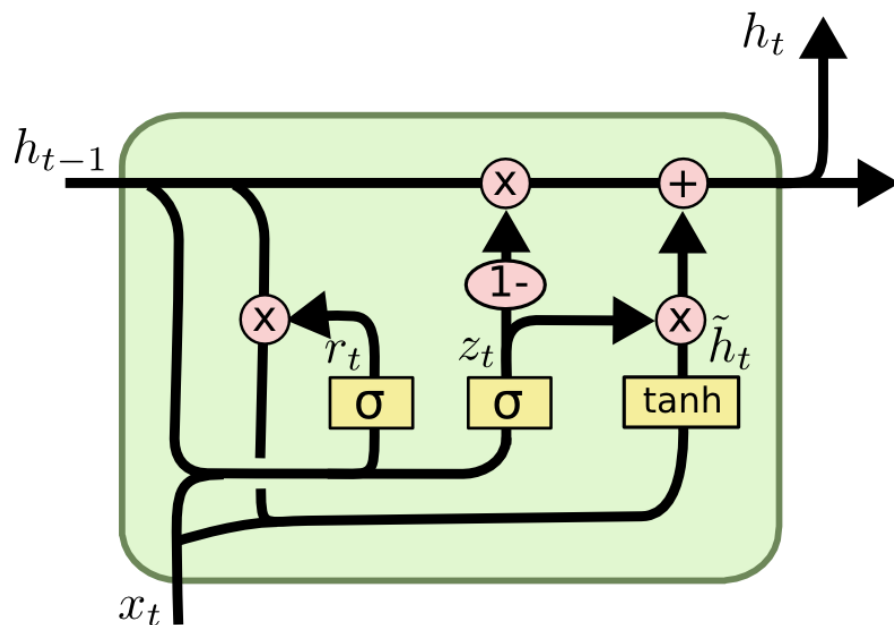


$$\mathbf{r}_t \leftarrow \sigma(\mathbf{W}^r \mathbf{x}_t + \mathbf{V}^r \mathbf{h}_{t-1} + \mathbf{b}^r)$$

$$\mathbf{u}_t \leftarrow \sigma(\mathbf{W}^u \mathbf{x}_t + \mathbf{V}^u \mathbf{h}_{t-1} + \mathbf{b}^u)$$

$$\hat{\mathbf{h}}_t \leftarrow \tanh(\mathbf{W}^h \mathbf{x}_t + \mathbf{V}^h (\mathbf{r}_t \cdot \mathbf{h}_{t-1}) + \mathbf{b}^h)$$

$$\mathbf{h}_t \leftarrow \mathbf{u}_t \cdot \mathbf{h}_{t-1} + (1 - \mathbf{u}_t) \cdot \hat{\mathbf{h}}_t$$



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

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

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

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

<http://colah.github.io/posts/2015-08-Understanding-LSTMs/img/LSTM3-var-GRU.png>

Highway Networks

For input \mathbf{x} , fully connected layer computes

$$\mathbf{y} \leftarrow H(\mathbf{x}, \mathbf{W}_H).$$

Highway networks add residual connection with gating:

$$\mathbf{y} \leftarrow H(\mathbf{x}, \mathbf{W}_H) \cdot T(\mathbf{x}, \mathbf{W}_T) + \mathbf{x} \cdot (1 - T(\mathbf{x}, \mathbf{W}_T)).$$

Usually, the gating is defined as

$$T(\mathbf{x}, \mathbf{W}_T) \leftarrow \sigma(\mathbf{W}_T \mathbf{x} + \mathbf{b}_T).$$

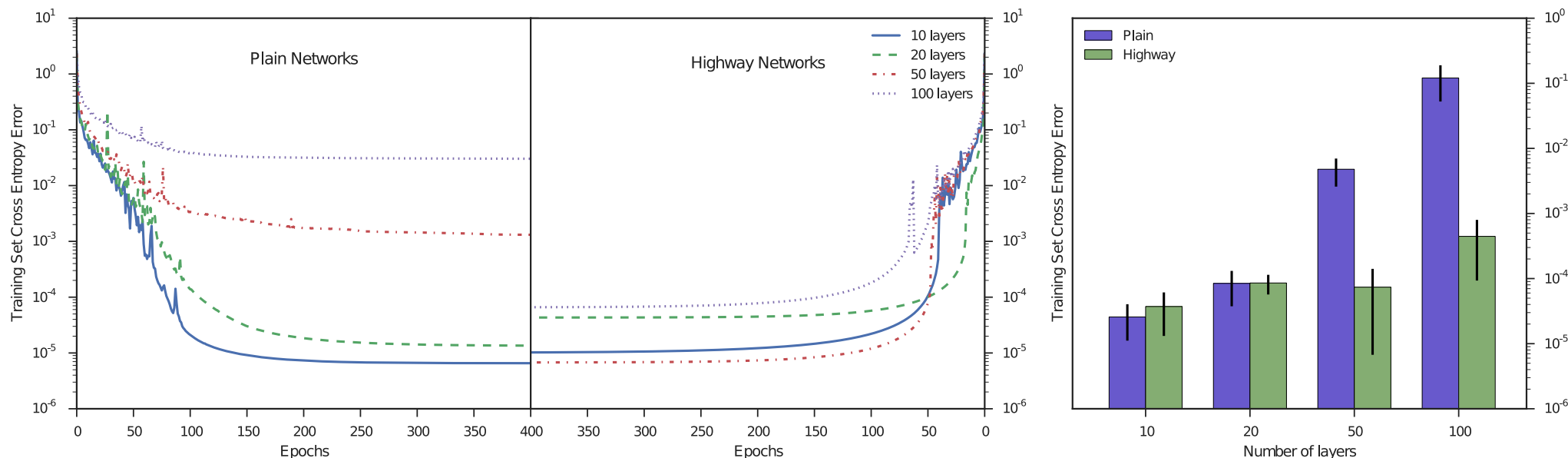


Figure 1: Comparison of optimization of plain networks and highway networks of various depths. *Left:* The training curves for the best hyperparameter settings obtained for each network depth. *Right:* Mean performance of top 10 (out of 100) hyperparameter settings. Plain networks become much harder to optimize with increasing depth, while highway networks with up to 100 layers can still be optimized well. Best viewed on screen (larger version included in Supplementary Material).

Figure 1 of paper "Training Very Deep Networks", <https://arxiv.org/abs/1507.06228>.

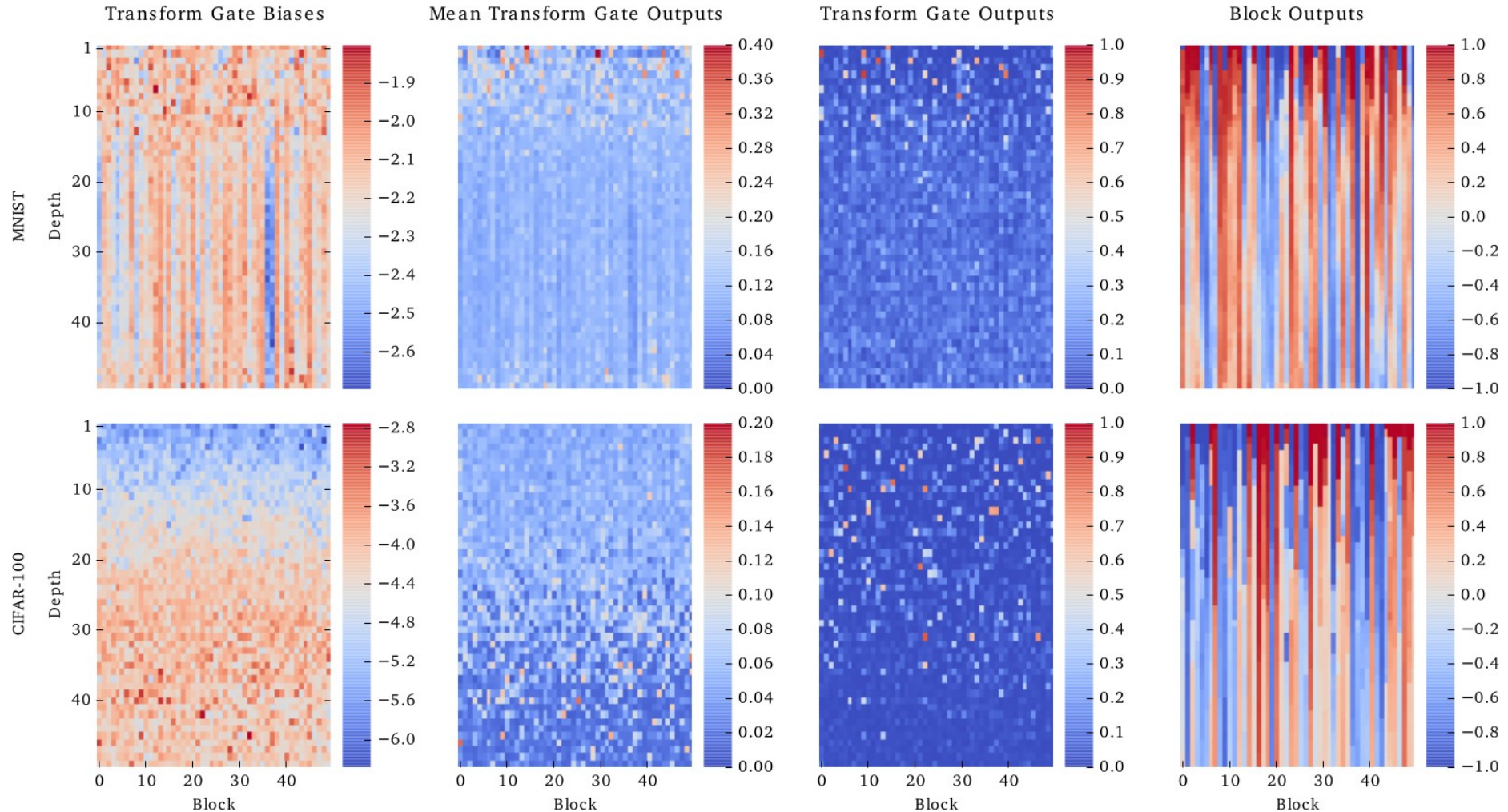


Figure 2 of paper "Training Very Deep Networks", <https://arxiv.org/abs/1507.06228>.

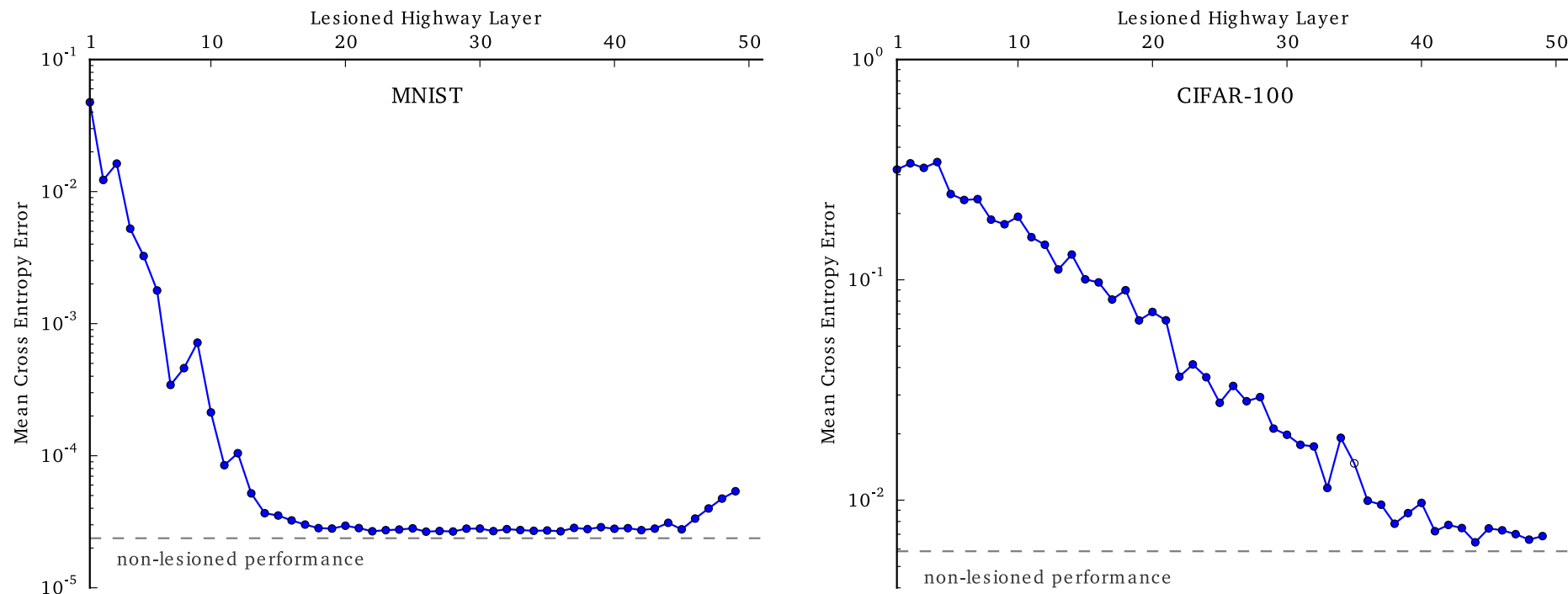
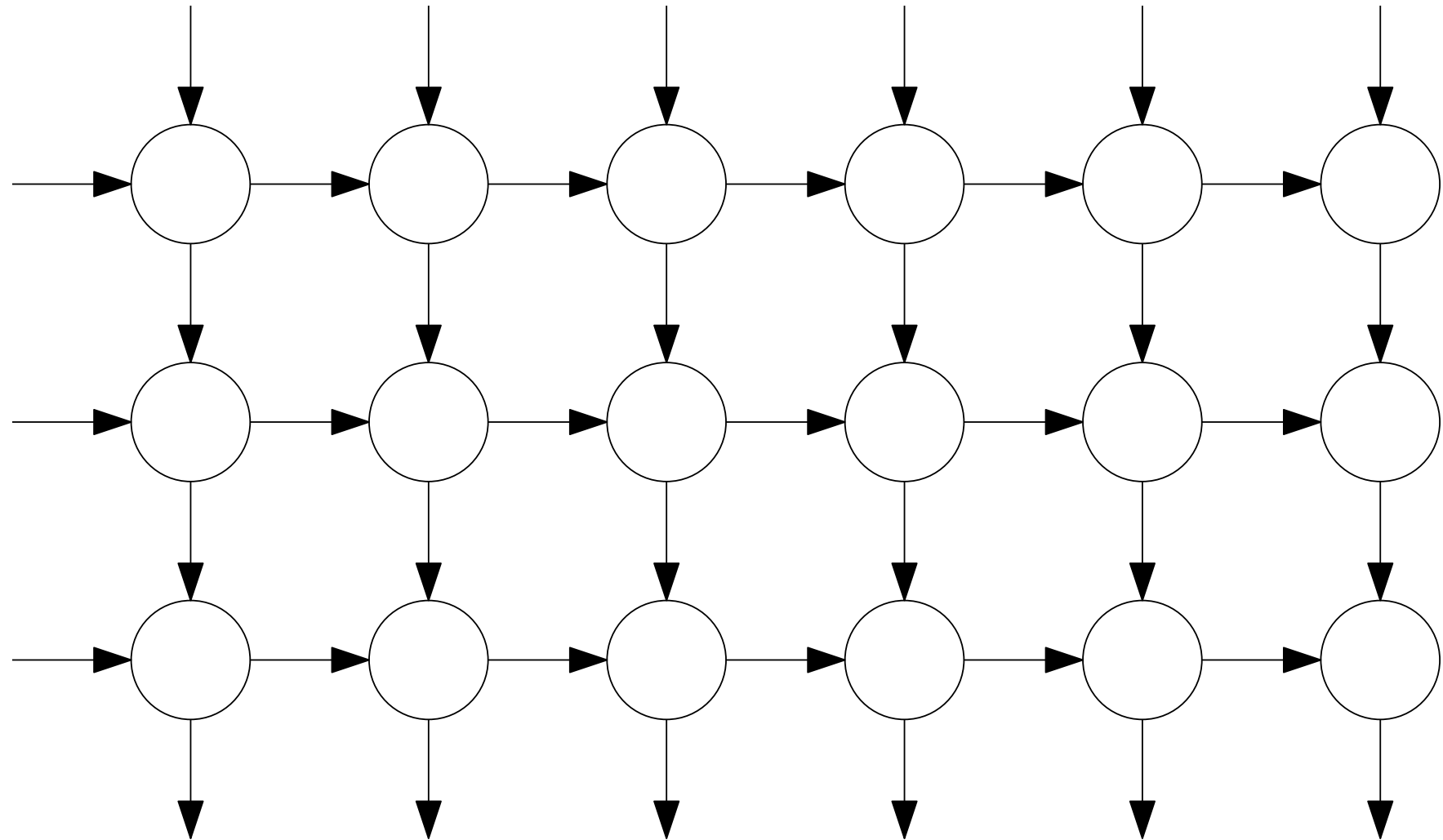
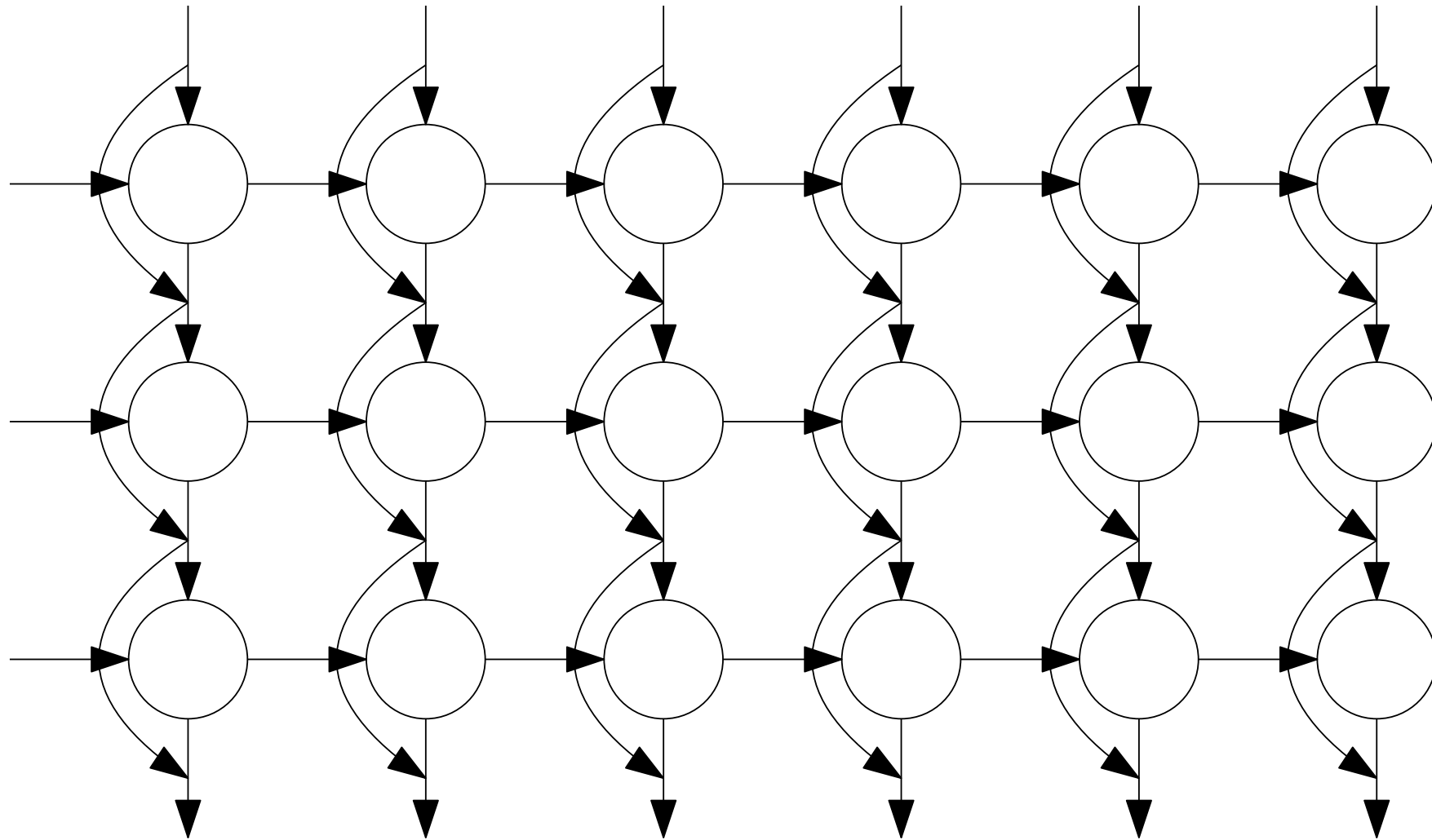


Figure 4: Lesioned training set performance (y-axis) of the best 50-layer highway networks on MNIST (left) and CIFAR-100 (right), as a function of the lesioned layer (x-axis). Evaluated on the full training set while forcefully closing all the transform gates of a single layer at a time. The non-lesioned performance is indicated as a dashed line at the bottom.

Figure 4 of paper "Training Very Deep Networks", <https://arxiv.org/abs/1507.06228>.

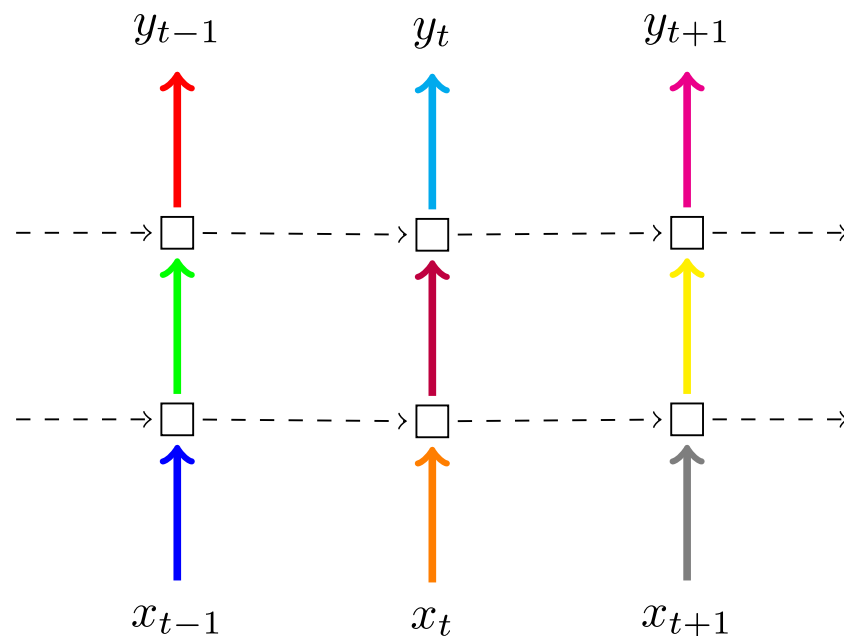




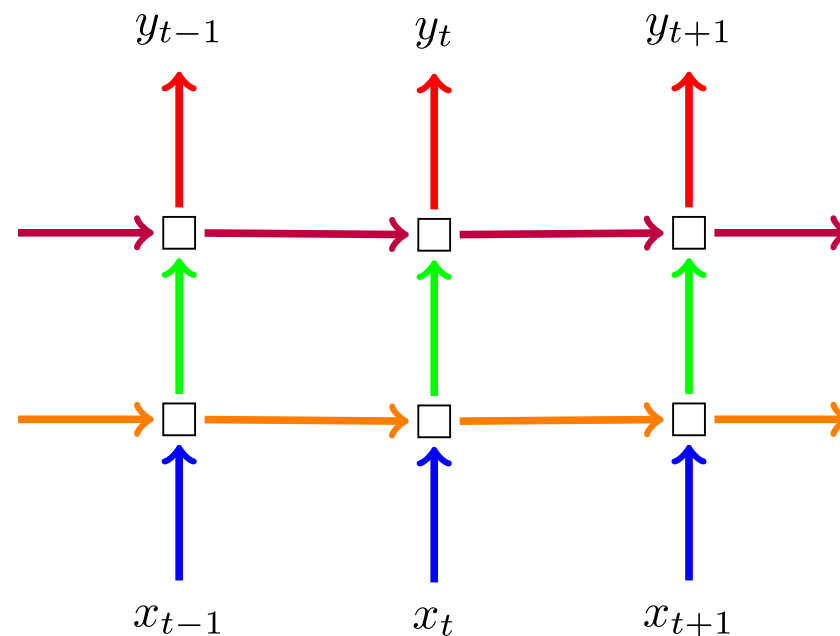
Dropout

- Using dropout on hidden states interferes with long-term dependencies.
- However, using dropout on the inputs and outputs works well and is used frequently.
 - In case residual connections are present, the output dropout needs to be applied before adding the residual connection.
- Several techniques were designed to allow using dropout on hidden states.
 - Variational Dropout
 - Recurrent Dropout
 - Zoneout

Variational Dropout



(a) Naive dropout RNN



(b) Variational RNN

Figure 1 of paper "A Theoretically Grounded Application of Dropout in Recurrent Neural Networks", <https://arxiv.org/abs/1512.05287.pdf>

Implemented in `tf.keras.layers.{RNN,LSTM,GRU}` using `dropout` and `recurrent_dropout` arguments (for dropping inputs and previous states, respectively).

Recurrent Dropout

Dropout only candidate states (i.e., values added to the memory cell in LSTM and previous state in GRU).

Zoneout

Randomly preserve hidden activations instead of dropping them.

Batch Normalization

Very fragile and sensitive to proper initialization (there were papers with negative results until people managed to make it work).

Layer Normalization

Much more stable than batch normalization.

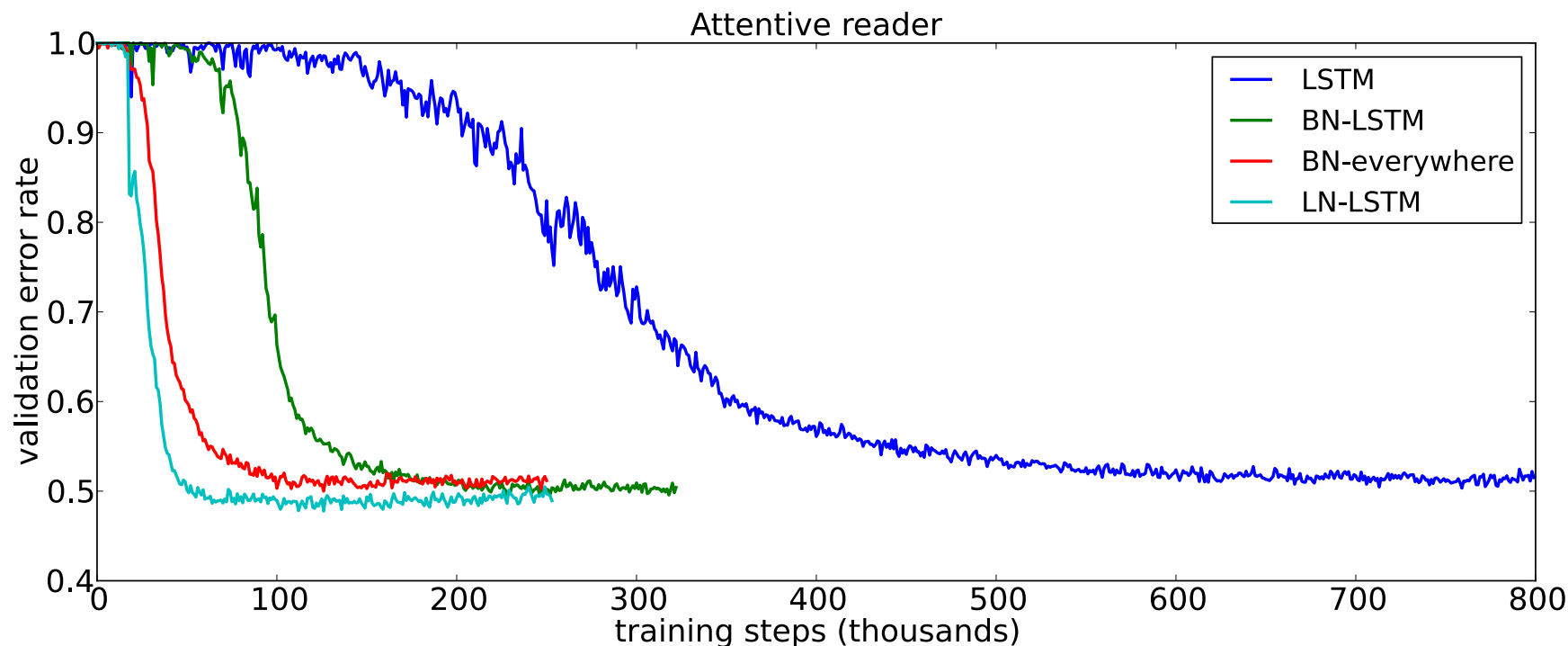


Figure 2: Validation curves for the attentive reader model. BN results are taken from [Cooijmans et al., 2016].

Figure 2 of paper "Layer Normalization", <https://arxiv.org/abs/1607.06450>.

One-hot encoding considers all words to be independent of each other.

However, words are not independent – some are more similar than others.

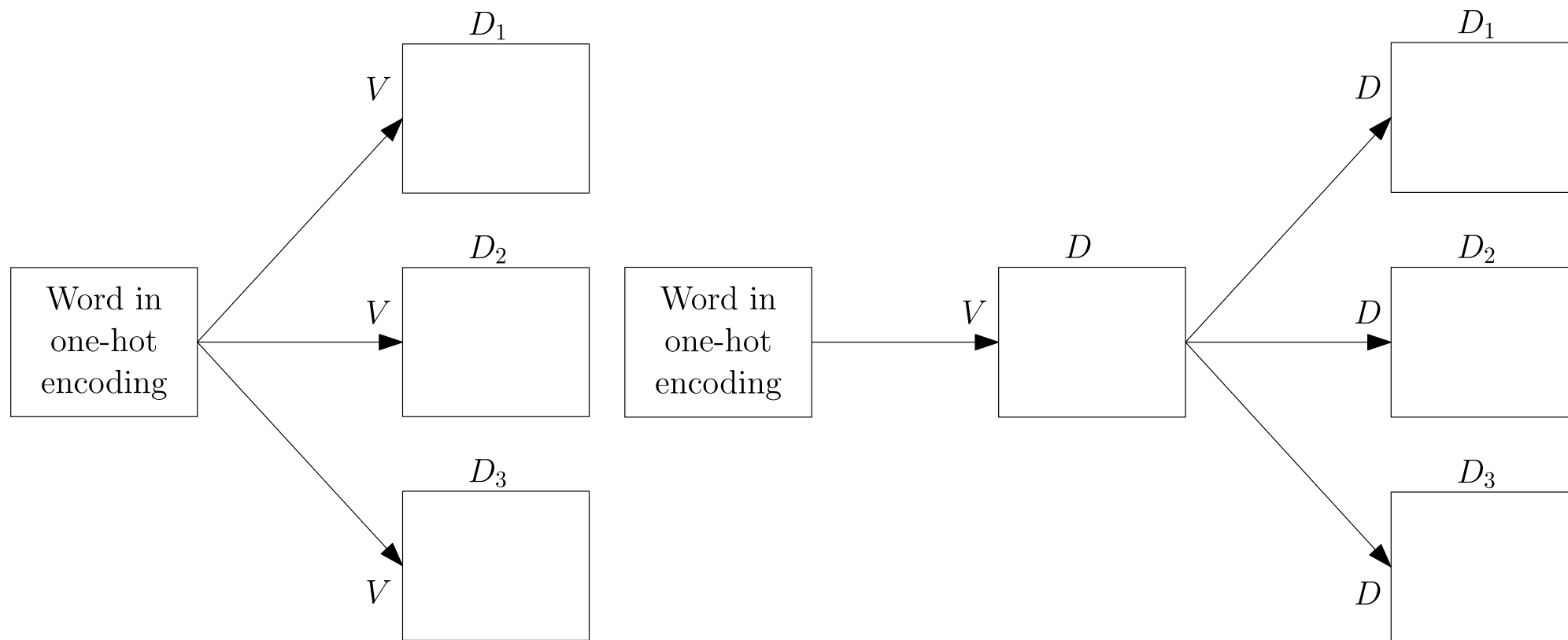
Ideally, we would like some kind of similarity in the space of the word representations.

Distributed Representation

The idea behind distributed representation is that objects can be represented using a set of common underlying factors.

We therefore represent words as fixed-size *embeddings* into \mathbb{R}^d space, with the vector elements playing role of the common underlying factors.

The word embedding layer is in fact just a fully connected layer on top of one-hot encoding. However, it is important that this layer is *shared* across the whole network.



Recurrent Character-level WEs

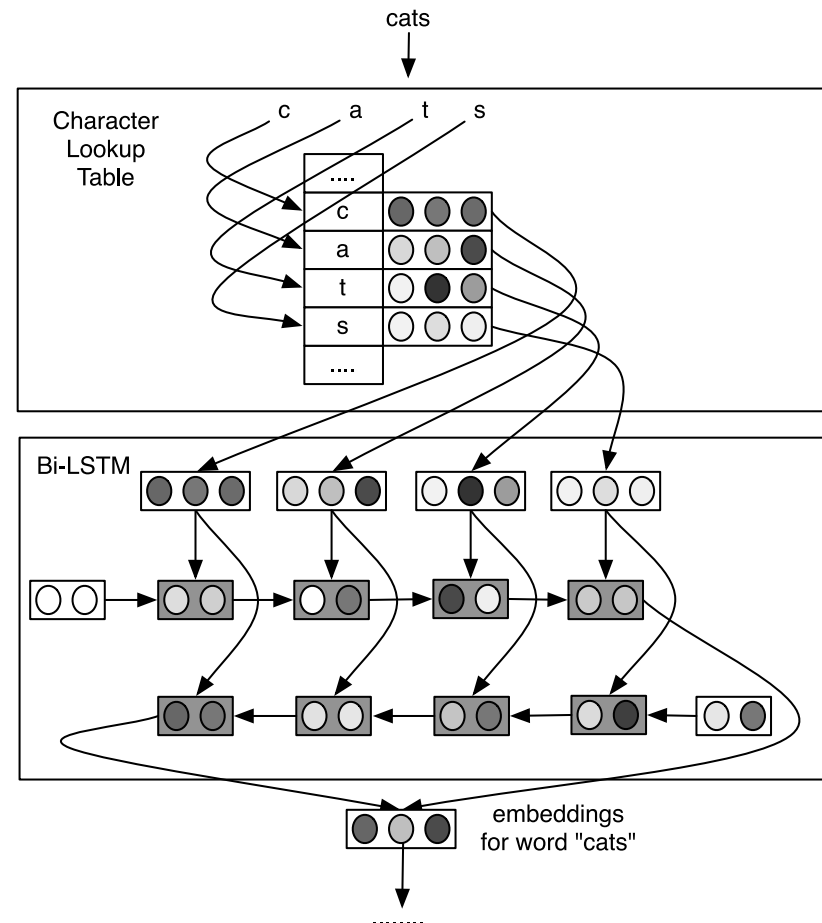


Figure 1 of paper "Finding Function in Form: Compositional Character Models for Open Vocabulary Word Representation", <https://arxiv.org/abs/1508.02096>.

Convolutional Character-level WEs

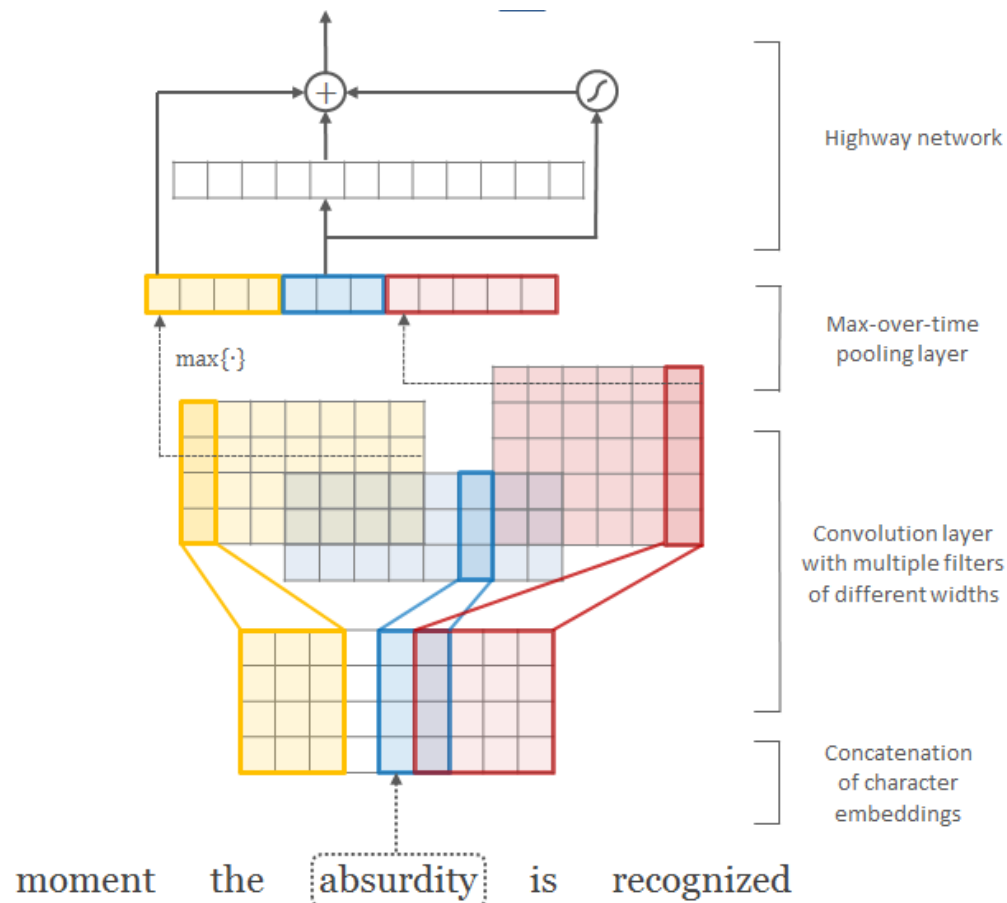


Figure 1 of paper "Character-Aware Neural Language Models", <https://arxiv.org/abs/1508.06615>.

Training

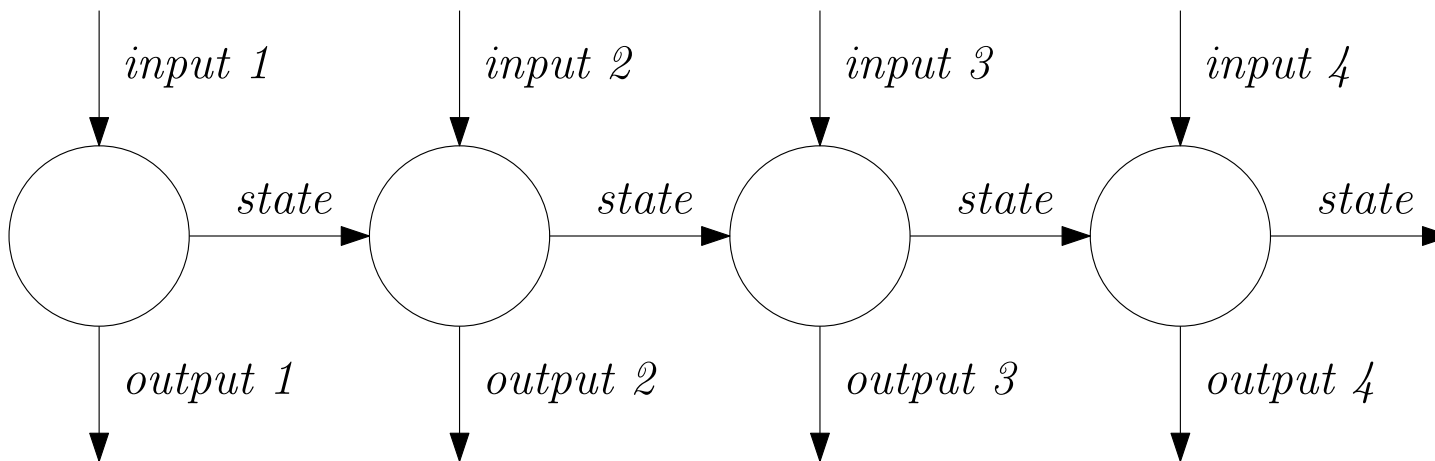
- Generate unique words per batch.
- Process the unique words in the batch.
- Copy the resulting embeddings suitably in the batch.

Inference

- We can cache character-level word embeddings during inference.

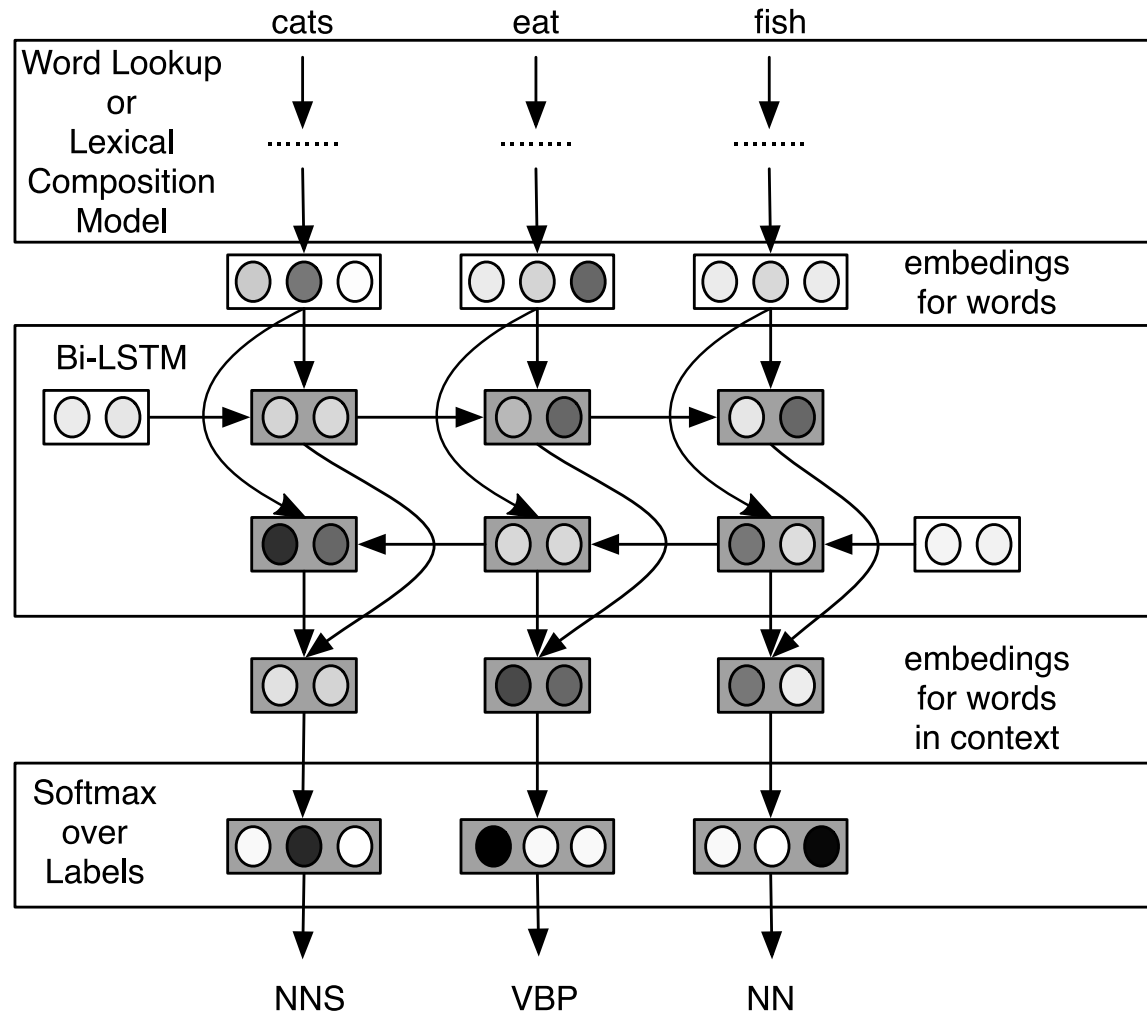
Sequence Element Classification

Use outputs for individual elements.



Sequence Representation

Use state after processing the whole sequence (alternatively, take output of the last element).



Modified Figure 3 of paper "Finding Function in Form: Compositional Character Models for Open Vocabulary Word Representation", <https://arxiv.org/abs/1508.02096>.

Sequence Tagging

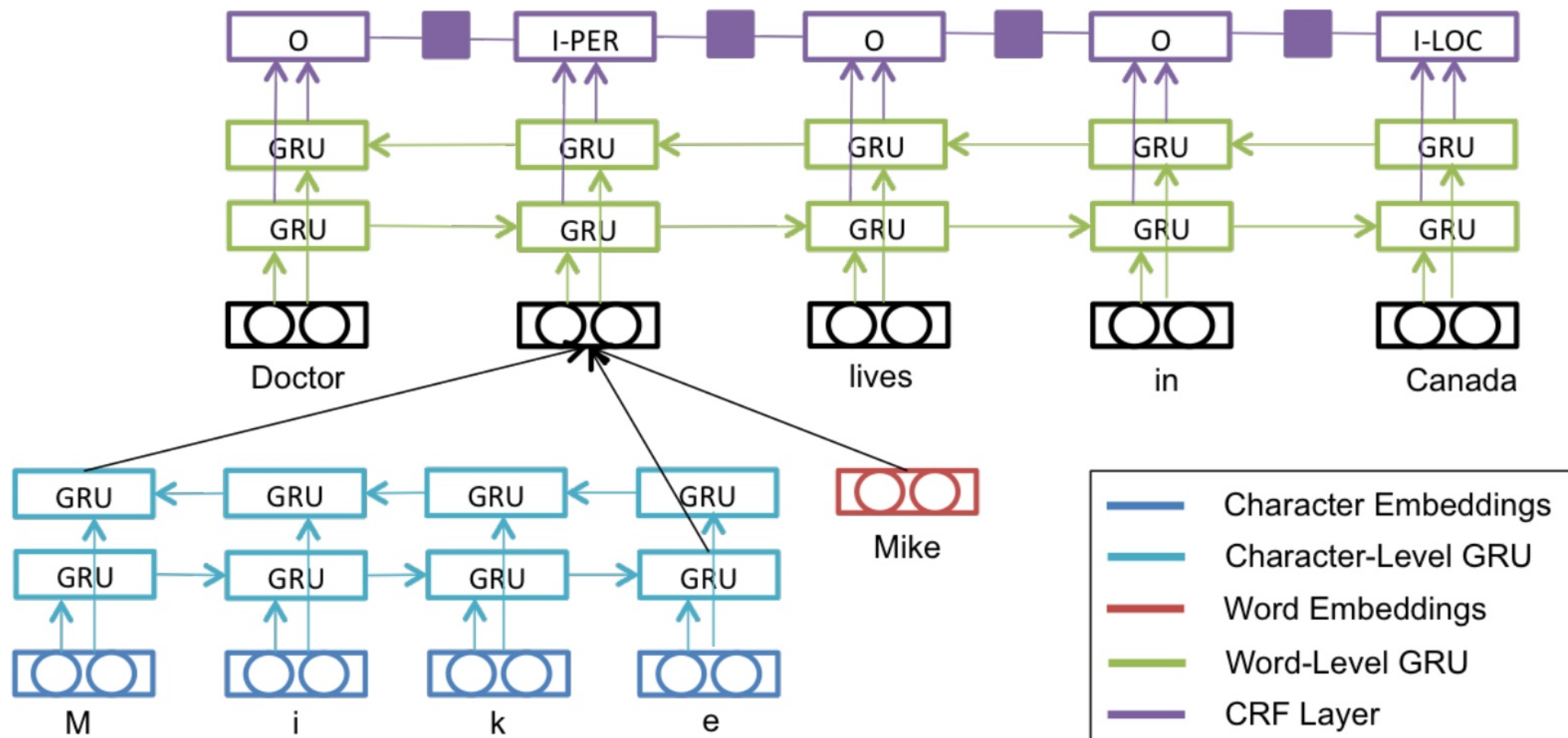
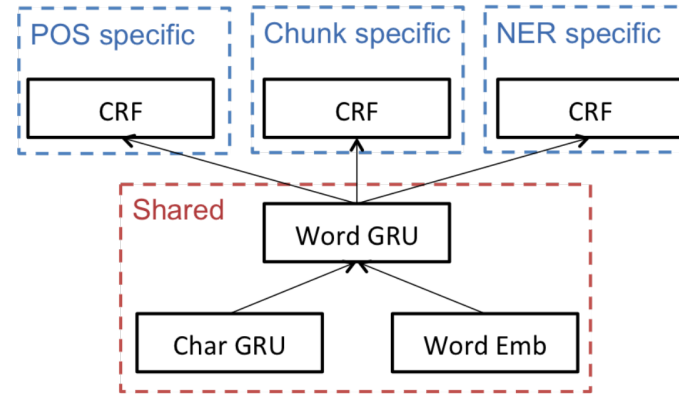
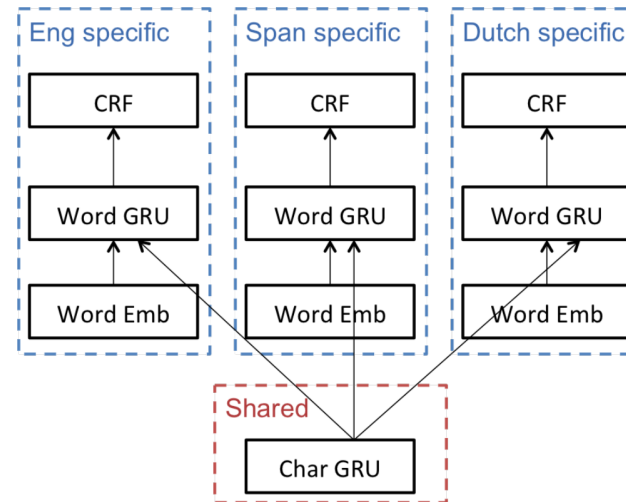


Figure 1 of paper "Multi-Task Cross-Lingual Sequence Tagging from Scratch", <https://arxiv.org/abs/1603.06270>.



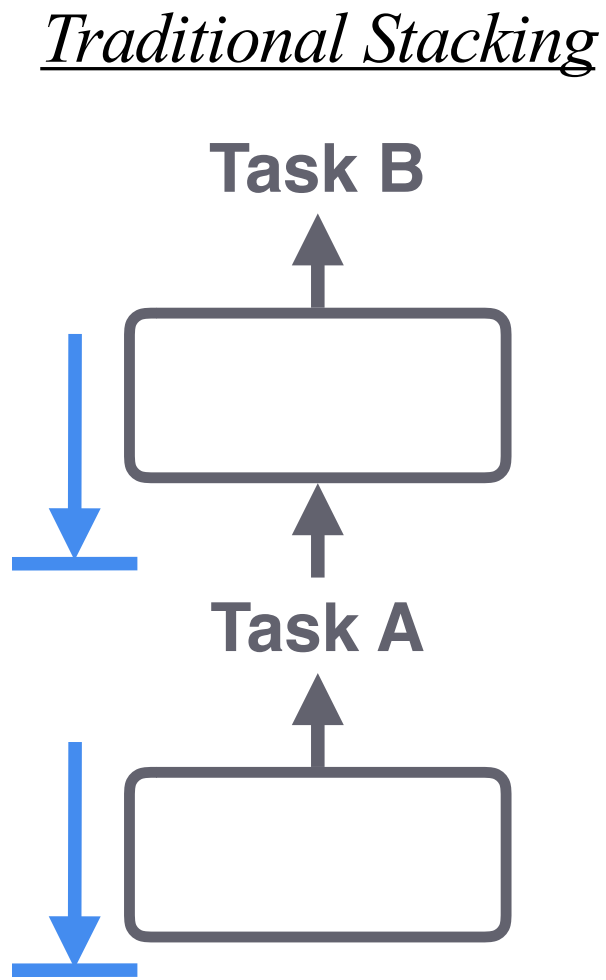
(a) Multi-Task Joint Training



(b) Cross-Lingual Joint Training

Figure 2 of paper "Multi-Task Cross-Lingual Sequence Tagging from Scratch", <https://arxiv.org/abs/1603.06270>.

Backpropagation



Stack-propagation

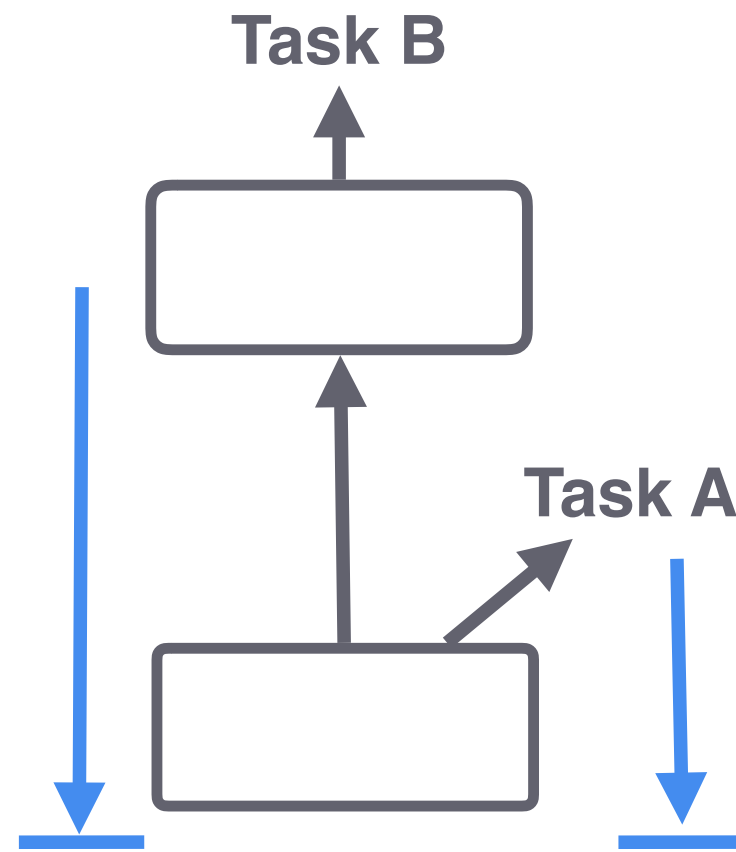


Figure 1 of paper "Stack-propagation: Improved Representation Learning for Syntax", <https://arxiv.org/abs/1603.06598>.

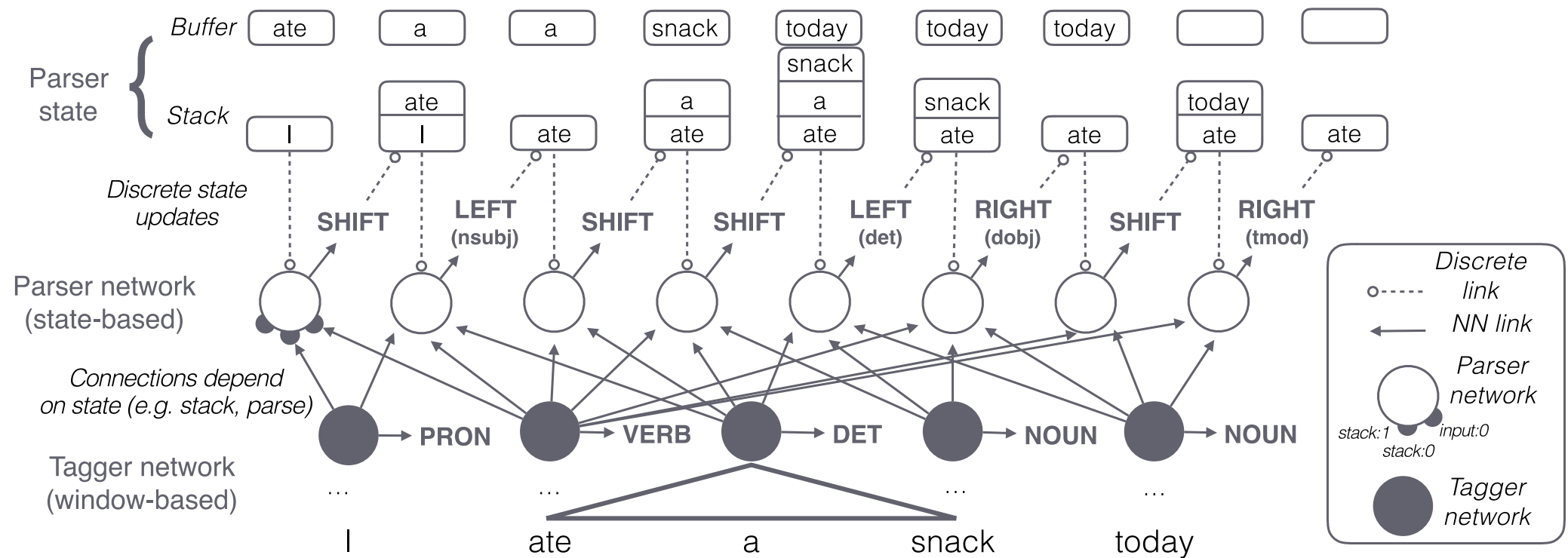


Figure 2 of paper "Stack-propagation: Improved Representation Learning for Syntax", <https://arxiv.org/abs/1603.06598>.

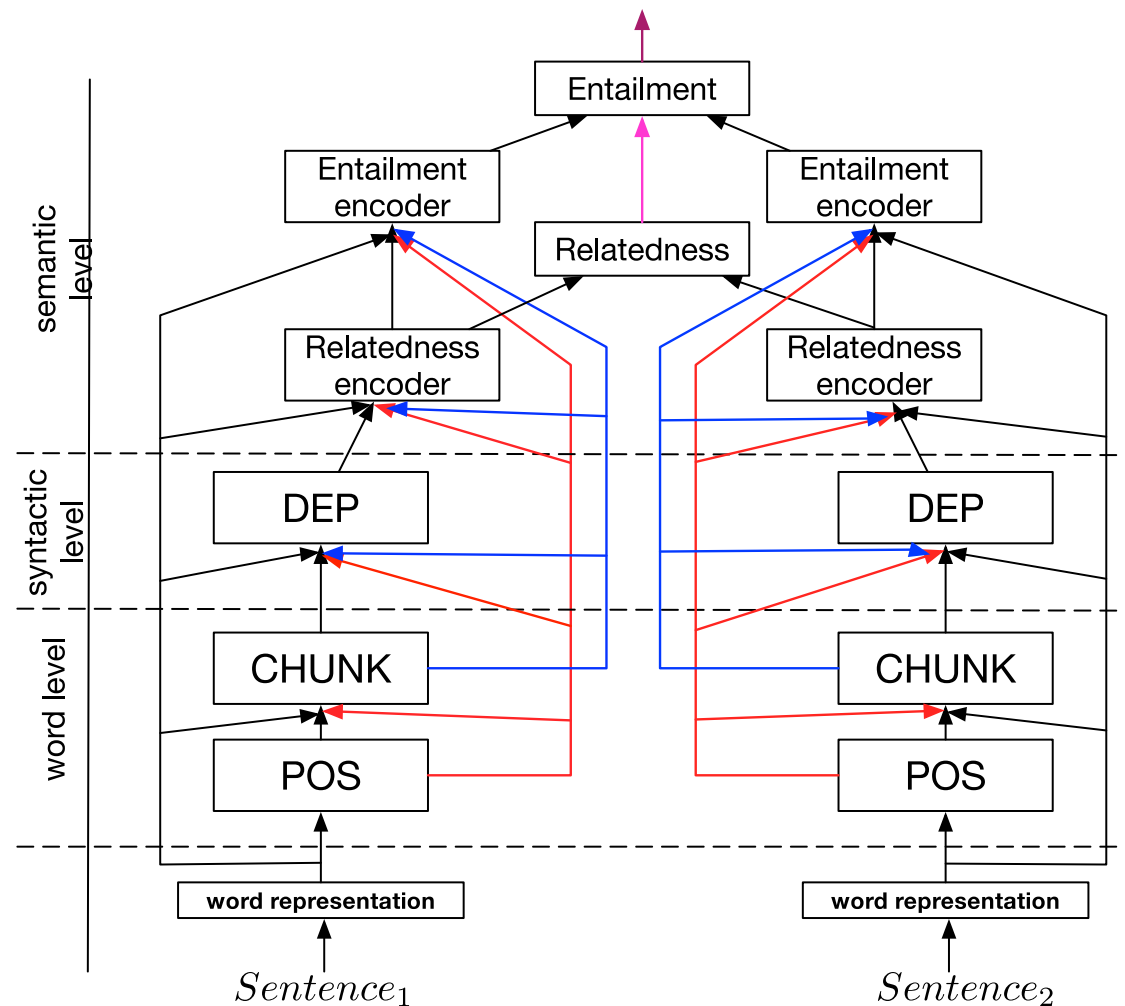


Figure 1 of paper "A Joint Many-Task Model: Growing a Neural Network for Multiple NLP Tasks", <https://arxiv.org/abs/1611.01587>.

Structured Prediction

Consider generating a sequence of $y_1, \dots, y_N \in Y^N$ given input $\mathbf{x}_1, \dots, \mathbf{x}_N$.

Predicting each sequence element independently models the distribution $P(y_i | \mathbf{X})$.

However, there may be dependencies among the y_i themselves, which is difficult to capture by independent element classification.

Let $G = (V, E)$ be a graph such that Y is indexed by vertices of G . Then (\mathbf{X}, \mathbf{y}) is a conditional Markov field, if the random variables \mathbf{y} conditioned on \mathbf{X} obey the Markov property with respect to the graph, i.e.,

$$P(y_i | \mathbf{X}, y_j, i \neq j) = P(y_i | \mathbf{X}, y_j \forall j : (i, j) \in E).$$

Usually we assume that dependencies of \mathbf{y} , conditioned on \mathbf{X} , form a chain.

Linear-chain Conditional Random Fields, usually abbreviated only to CRF, acts as an output layer. It can be considered an extension of a softmax – instead of a sequence of independent softmaxes, CRF is a sentence-level softmax, with additional weights for neighboring sequence elements.

$$s(\mathbf{X}, \mathbf{y}; \boldsymbol{\theta}, \mathbf{A}) = \sum_{i=1}^N (\mathbf{A}_{y_{i-1}, y_i} + f_{\boldsymbol{\theta}}(y_i | \mathbf{X}))$$

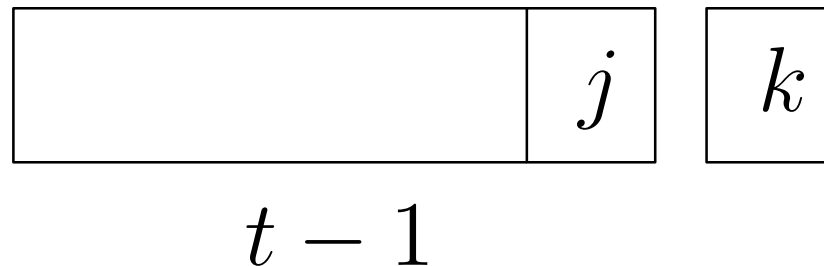
$$p(\mathbf{y} | \mathbf{X}) = \text{softmax}_{\mathbf{z} \in Y^N} (s(\mathbf{X}, \mathbf{z}))_{\mathbf{z}}$$

$$\log p(\mathbf{y} | \mathbf{X}) = s(\mathbf{X}, \mathbf{y}) - \text{logadd}_{\mathbf{z} \in Y^N} (s(\mathbf{X}, \mathbf{z}))$$

Computation

We can compute $p(\mathbf{y}|\mathbf{X})$ efficiently using dynamic programming. We denote $\alpha_t(k)$ the logarithmic probability of all t -element sequences with the last label y being k .

The core idea is the following:



$$\alpha_t(k) = f_{\theta}(y_t = k|\mathbf{X}) + \text{logadd}_{j \in Y}(\alpha_{t-1}(j) + \mathbf{A}_{j,k}).$$

For efficient implementation, we use the fact that

$$\ln(a + b) = \ln a + \ln(1 + e^{\ln b - \ln a}).$$

Inputs: Network computing $f_{\theta}(y_t = k | \mathbf{X})$, an unnormalized probability of output sequence element probability being k at time t .

Inputs: Transition matrix $\mathbf{A} \in \mathbb{R}^{Y \times Y}$.

Inputs: Input sequence \mathbf{X} of length N , gold labeling $\mathbf{y}^g \in Y^N$.

Outputs: Value of $\log p(\mathbf{y} | \mathbf{X})$.

Time Complexity: $\mathcal{O}(N \cdot Y^2)$.

- For $t = 1, \dots, N$:
 - For $k = 1, \dots, Y$:
 - $\alpha_t(k) \leftarrow f_{\theta}(y_t = k | \mathbf{X})$
 - If $t > 1$:
 - For $j = 1, \dots, Y$:
 - $\alpha_t(k) \leftarrow \text{logadd}(\alpha_t(k), \alpha_{t-1}(j) + \mathbf{A}_{j,k})$
- Return $\sum_{t=1}^N f_{\theta}(y_t = y_t^g | \mathbf{X}) + \sum_{t=2}^N \mathbf{A}_{y_{t-1}^g, y_t^g} - \text{logadd}_{k=1}^Y(\alpha_N(k))$

Decoding

We can perform optimal decoding, by using the same algorithm, only replacing `logadd` with `max` and tracking where the maximum was attained.

Applications

CRF output layers are useful for *span labeling* tasks, like

- named entity recognition
- dialog slot filling