# Recurrent Neural Networks

**Milan Straka**

📅 **April 8, 2024**

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

# Recurrent Neural Networks

*– koncipování pro schvence, jejichž jedna dimenze je čas*

# Recurrent Neural Networks

## Single RNN cell

*input*

*state*

*output*

držím si nějaký stav v tí vyhodnocovací
buňce

## Unrolled RNN cells

*input 1*     *input 2*     *input 3*     *input 4*

*state*     *state*     *state*     *state*

*output 1*     *output 2*     *output 3*     *output 4*

ten state si předávám

tohle mi má reprezentovat
historii toho vstupa

*Použilo se už před 30 lety*

*input*

→ *mohl byoh udělat fully-connected vrstvu*

*previous state*

*output = new state*

*input*

tanh *state*

*output*

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

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

One of the simplest possibilities (called `SimpleRNN` in Keras, `RNN` in PyTorch) is

*jasto to ale bude blbě propagovat zpětně nějakou chybu*

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

*tanh puto, že je omezený. ReLU je neomezené...*

# Basic RNN Cell

Basic RNN cells suffer a lot from vanishing/exploding gradients (the so-called **challenge of long-term dependencies**). *tohle už je problém    tanh tohle pojistí protože max(tanh(x)) = 1*

If we simplify the recurrence of states to just a linear approximation

$$h^{(t)} \approx U h^{(t-1)},$$

we get $h^{(t)} \approx U^t h^{(0)}$.

*→ tohle by znamenalo, že kdyby měla největší nějakou dimenzi potlačovat, tak na základě „t" by ji potlačila hrozně moc.*

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

$$h^{(t)} \approx Q \Lambda^t Q^{-1} h^{(0)}.$$

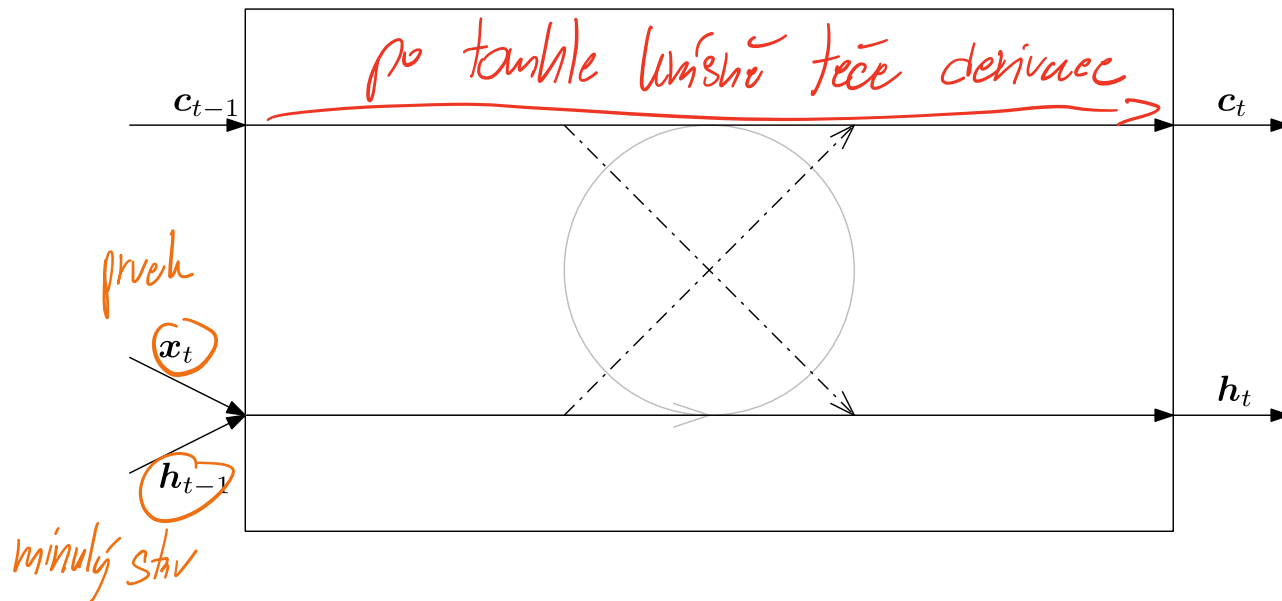The main problem is that the *same* function is iteratively applied many times.

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

Hochreiter & Schmidhuber (1997) suggested that to enforce *constant error flow*, we would like

$$f' = 1.$$

*kdyz ten gradient pujde zpatky, hh bude zustavat stejny*

They propose to achieve that by a *constant error carrousel*.

*po tamhle lanisho tece derivace*

*prveh*

$x_t$

$h_{t-1}$

*minulý stav*

$c_{t-1}$

$c_t$

$h_t$

# Long Short-Term Memory

They also propose an **input** and **output** gates which control the flow of information into and out of the carrousel (**memory cell $c_t$**).
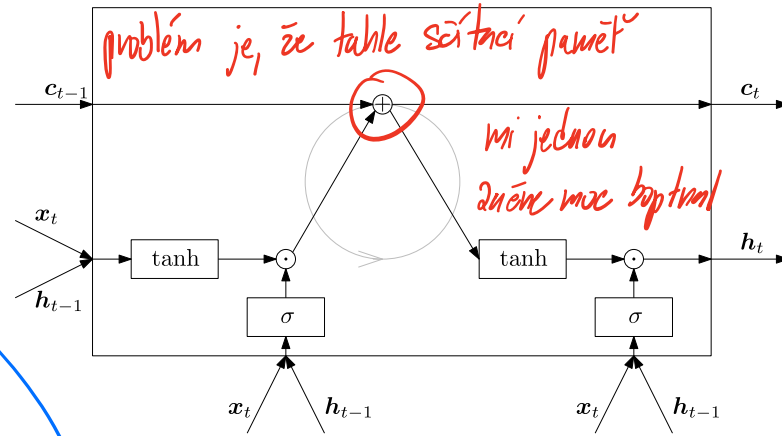
$$i_t \leftarrow \sigma(W^i x_t + V^i h_{t-1} + b^i)$$

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

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

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

*sigmoid*

*tahle je tu dřýv RNN cell.*

*problém je, že tahle sčítací paměť*

*mi jednou zrůne moc bogtnal*

*pro kuždou dimenzi si vytvořím Past,*

*jak moc se mi dané dato v dimenzi hodí*

*nastavnji, které hodnoty chci jak moc posílat ven.*

*tanh proto, že $c_t$ by jinak mohlo střílet vysoko...*



$c_{t-1}$     $c_t$

$x_t$

$h_{t-1}$      tanh    $\odot$    tanh    $\odot$    $h_t$

$\sigma$      $\sigma$

$x_t$   $h_{t-1}$      $x_t$   $h_{t-1}$

# Long Short-Term Memory

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

*hlavní je nezačít s 0 !*

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

*divence, kterou chci zapomenout*

*→ s tímhle se musí opatrně pracvat →*

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

*→ proto bias je dobrý = 1*

*brání se efektivně učit*

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

*hlavně na začátku*

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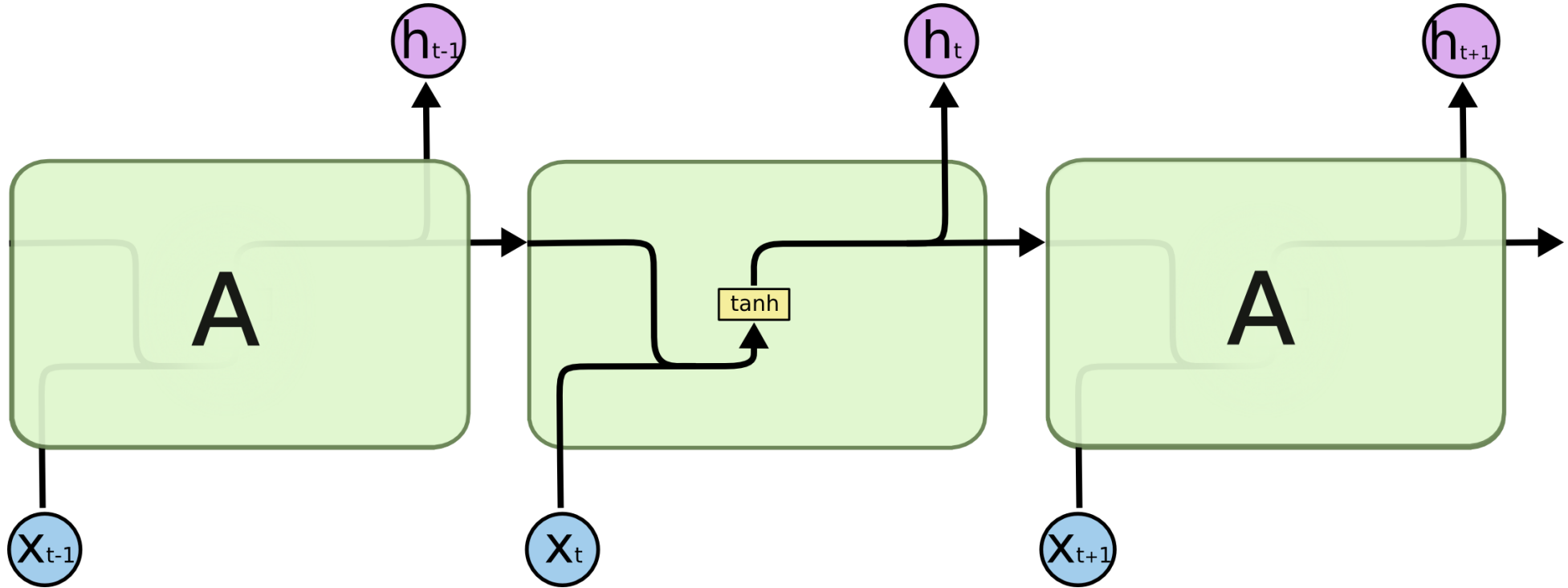
$$h_t \leftarrow o_t \odot \tanh(c_t)$$

Note that since 2015, following the paper

*Tohle je nynější LSTM*

- R. Jozefowicz et al.: *An Empirical Exploration of Recurrent Network Architectures*

the forget gate bias $b^f$ is usually initialized to 1, so that the forget gate is closer to 1 and the gradients can easily flow through multiple timesteps. (Gers et al. advocated this in the original paper already.) (BTW, I think 3 might be even better, as $\sigma(1) \approx 0.731, \sigma(3) \approx 0.953$.)
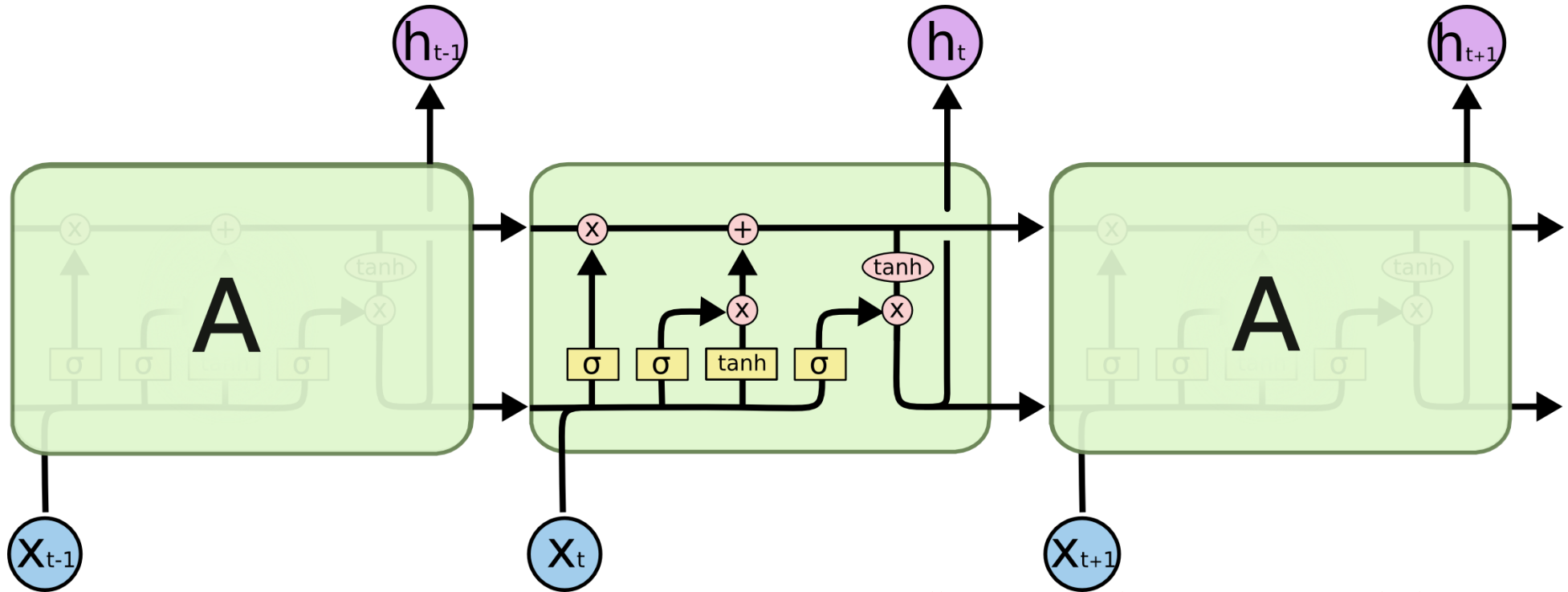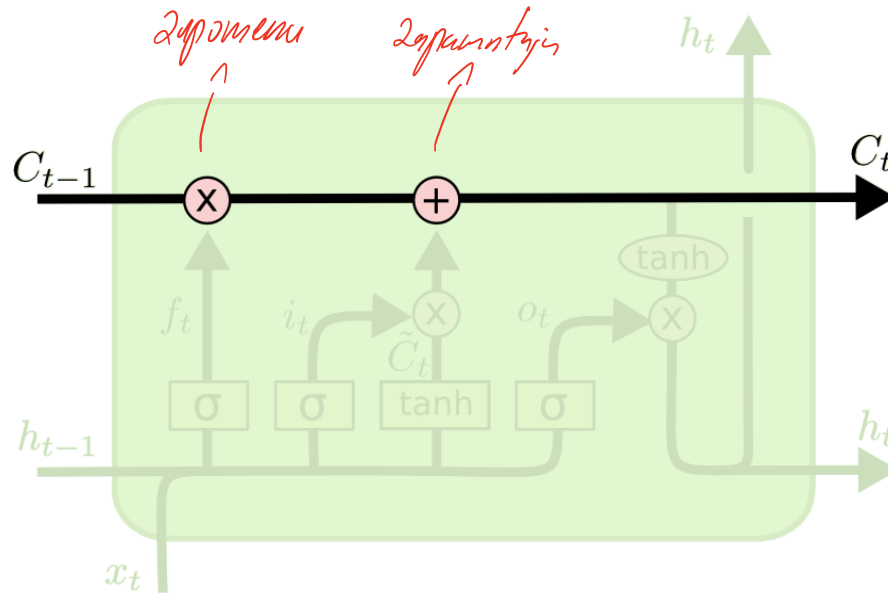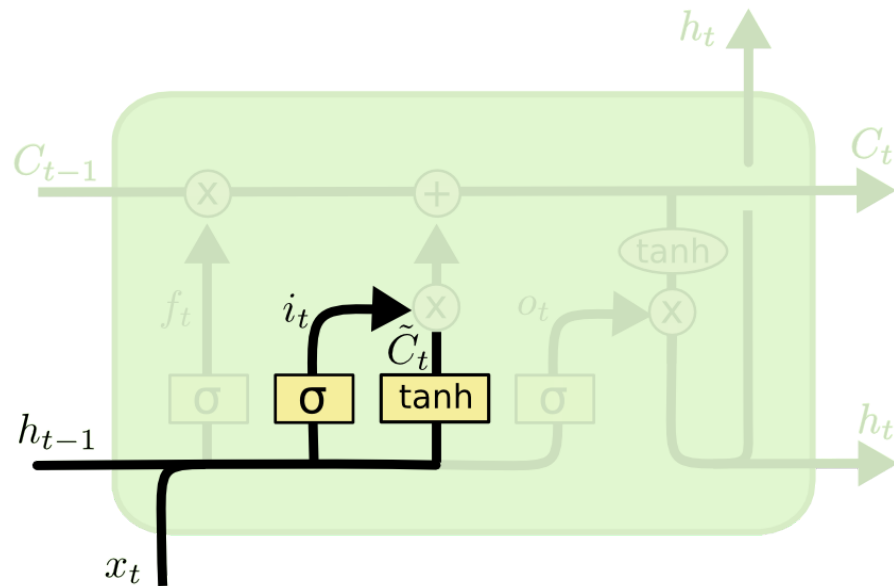
Simple RNN

*LSTM*



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

tohle je pouze paměťová buňka, která občas něco zapomene, jak prostě tlačí informaci dál...



zapomenu

zapamatuji

jak moc silně si informace chci akumulovat

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

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

kandidátní hodnota na akumulování

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

tohle můžu chápat jako
fully-connected vrstva

# Long Short-Term Memory



tohle mi určí, kterou část chci zapomenout.

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

# Long Short-Term Memory



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

*zapomenu*

*naučím se*

Tohle je z roku 2000 a dotud se nemělo nic dalšího lepšího.



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

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

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

Je tam ale hodně vah a proto se to dlouho trémuje

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

*forget je doplněk learnu*

- no memory cell,
- forgetting and updating tied together.

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

*čím tu chci nahradit*

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

*co chci nahradit / kde se chci učit*

$$\hat{h}_t \leftarrow \tanh(W^h x_t + V^h(r_t \odot h_{t-1}) + b^h)$$

$$h_t \leftarrow u_t \odot h_{t-1} + (1 - u_t) \odot \hat{h}_t$$

*vezmu si jen ḱus informace z minula*

*updatnu celý view, čast bude nahmcení*

*Fungyji ḱépr jak LSTM díkg umíě vahím*

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

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

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

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

doplňky

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

# GRU and LSTM Differences

The main differences between GRU and LSTM:

- GRU uses fewer parameters and less computation.
  - six matrices $W$, $V$ instead of eight

- GRU are easier to work with, because the state is just one tensor, while it is a pair of tensors for LSTM.

- In most tasks, LSTM and GRU give very similar results.

- However, there are some tasks, on which LSTM achieves (much) better results than GRU.
  - For a demonstration of difference in the expressive power of LSTM and GRU (caused by the coupling of the forget and update gate), see the paper
    - G. Weiss et al.: *On the Practical Computational Power of Finite Precision RNNs for Language Recognition* https://arxiv.org/abs/1805.04908

  - For a difference between LSTM and GRU on a real-word task, see for example
    - T. Dozat et al.: *Deep Biaffine Attention for Neural Dependency Parsing* https://arxiv.org/abs/1611.01734

GRU se neumí naučit nový věci, aniž by zapomněla nějaký stany

Recall that when we approximate $\boldsymbol{h}^{(t)} \approx \boldsymbol{U}\boldsymbol{h}^{(t-1)}$, assuming the eigenvalue decomposition of $\boldsymbol{U} = \boldsymbol{Q}\boldsymbol{\Lambda}\boldsymbol{Q}^{-1}$, we get

*to zajistí, že tahle mocnice bude jen otáčet data, nebude je ničem zvětšovat.*

$$\boldsymbol{h}^{(t)} \approx \boldsymbol{Q}\boxed{\boldsymbol{\Lambda}^t}\boldsymbol{Q}^{-1}\boldsymbol{h}^{(0)}.$$

This motivated a specific initialization scheme for the $\boldsymbol{U}$ matrix – this so-called **recurrent kernel** (the concatenation of all the $\boldsymbol{V}^i$, $\boldsymbol{V}^f$, $\boldsymbol{V}^o$, $\boldsymbol{V}^y$ matrices) is initialized with a randomly generated orthogonal matrix.

This **orthogonal** initialization is used for all RNN cells in Keras (via the `recurrent_initializer='orthogonal'` parameter of `SimpleRNN`, `GRU`, and `LSTM`).

# Highway Networks

# Highway Networks

For input $\boldsymbol{x}$, fully connected layer computes

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

Highway networks add residual connection with gating:

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

Usually, the gating is defined as         *nezhoduje, jak moc chci tu informaci tlačit dopředa*

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

Note that the resulting update is very similar to a GRU cell with $\boldsymbol{h}_t$ removed; for a fully connected layer $H(\boldsymbol{x}, \boldsymbol{W}_H) = \tanh(\boldsymbol{W}_H \boldsymbol{x} + \boldsymbol{b}_H)$ it is exactly it, apart from copying $\boldsymbol{x}$ instead of $\boldsymbol{h}_{t-1}$.

Analogously to LSTM, the transform gate bias $\boldsymbol{b}_T$ should be initialized to a negative number.
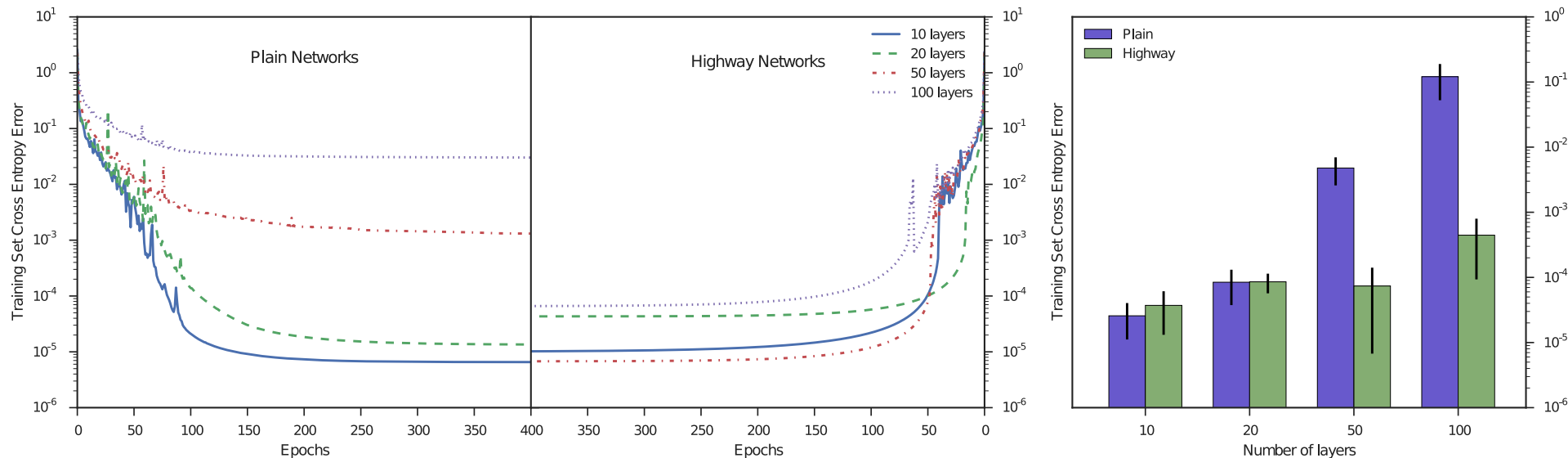
Figure 1: Comparison of optimization of plain networks and highway networks of various depths. *Left:* The training curves for the best hyperparameter settings obtained for each network depth. *Right:* Mean performance of top 10 (out of 100) hyperparameter settings. Plain networks become much harder 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).
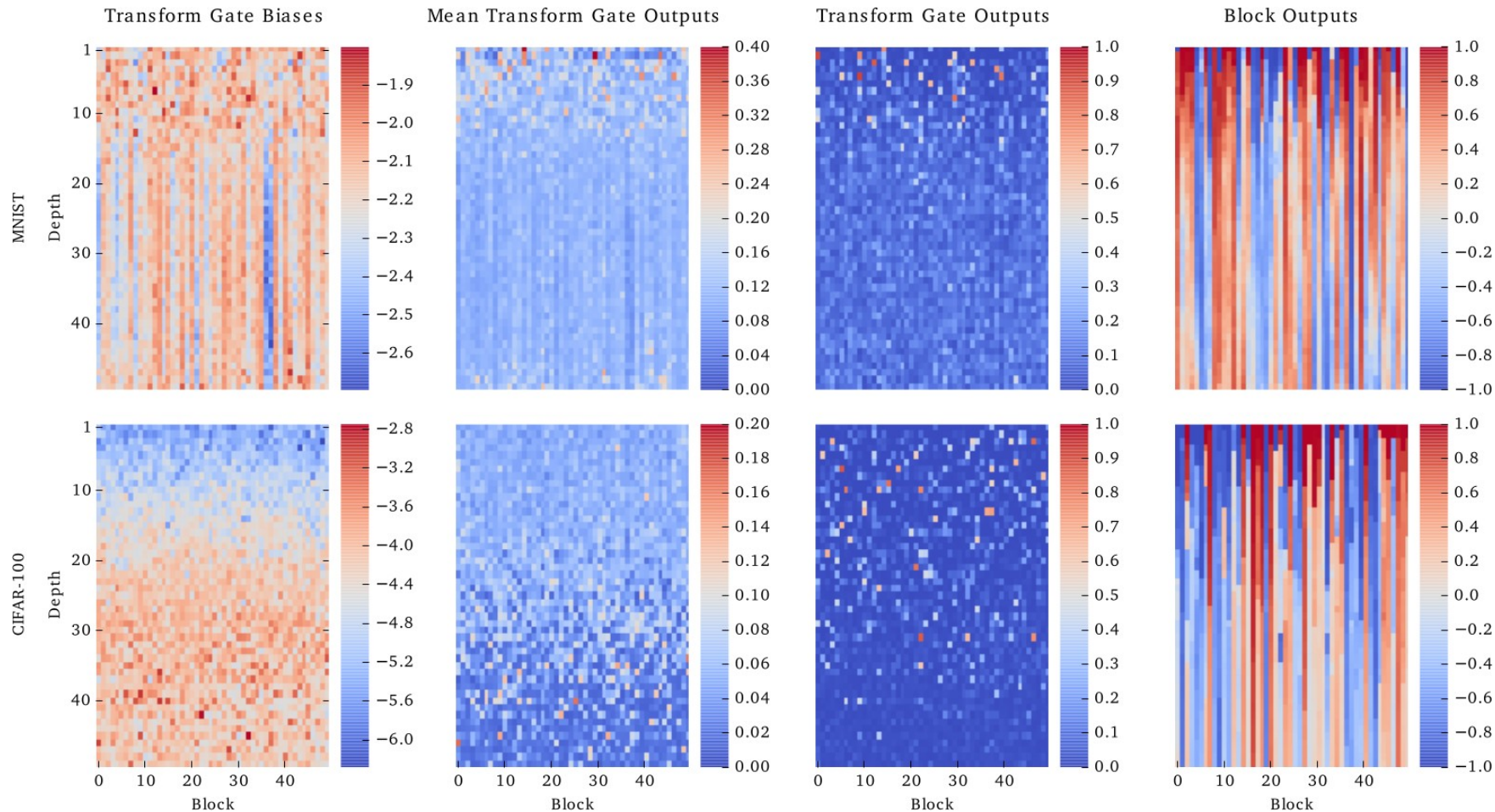
*Figure 1 of "Training Very Deep Networks", https://arxiv.org/abs/1507.06228*

# Highway Networks

Figure 2 of "Training Very Deep Networks", https://arxiv.org/abs/1507.06228
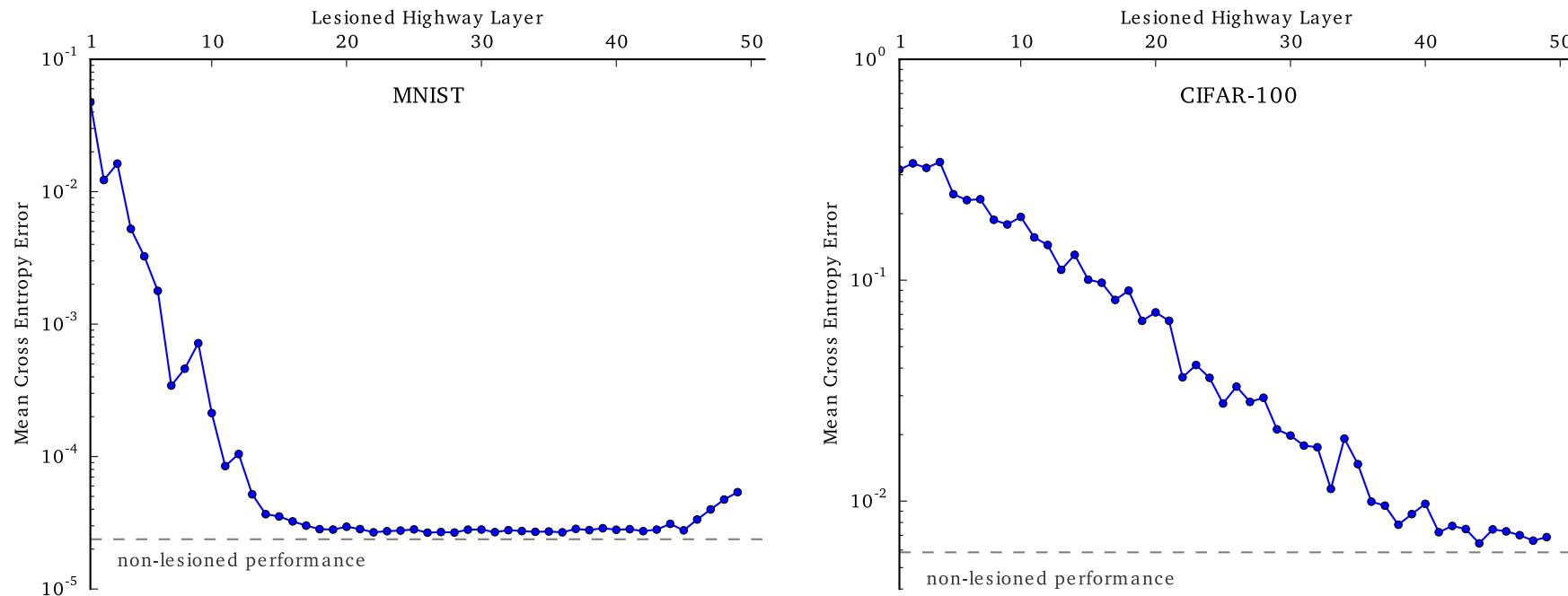
Figure 4: Lesioned training set performance (y-axis) of the best 50-layer highway networks on MNIST (left) and CIFAR-100 (right), as a function of the lesioned layer (x-axis). Evaluated on the full training set while forcefully closing all the transform gates of a single layer at a time. The non-lesioned performance is indicated as a dashed line at the bottom.

*Figure 4 of "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

nesmí se ale používat ve stavové informaci, protože bych efektivně všechny zapomněl, přestože jsem se to hrozně snažil naučit.
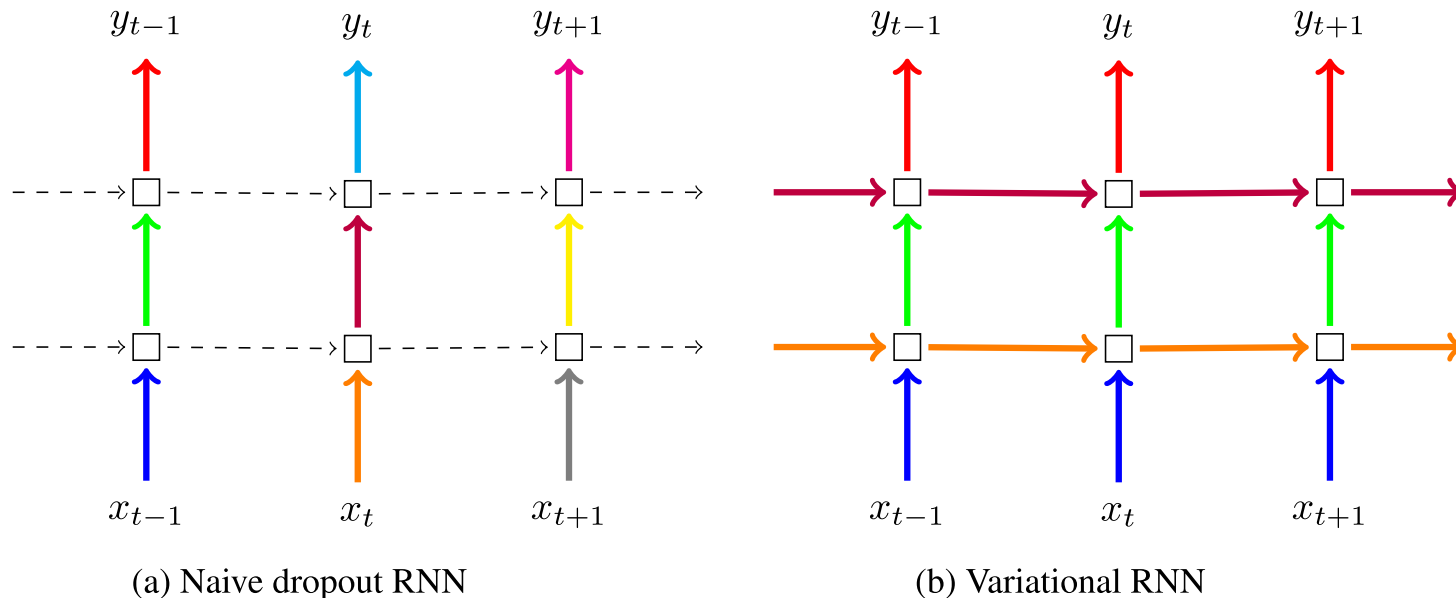
## Variational Dropout



(a) Naive dropout RNN          (b) Variational RNN

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

To implement variational dropout on inputs in Keras, use `noise_shape` of `keras.layers.Dropout` to force the same mask across time-steps. The variational dropout on the hidden states can be implemented using `recurrent_dropout` argument of `keras.layers.{LSTM,GRU,SimpleRNN}{,Cell}`.
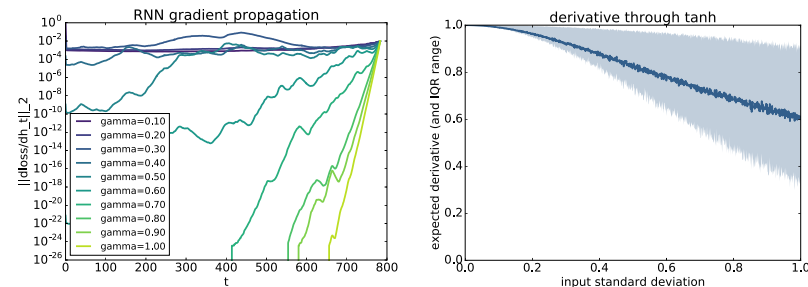
## Recurrent Dropout

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

## Zoneout

Randomly preserve hidden activations instead of dropping them.

## Batch Normalization

Very fragile and sensitive to proper initialization – there were papers with negative results (*Dario Amodei et al, 2015: Deep Speech 2* or *Cesar Laurent et al, 2016: Batch Normalized Recurrent Neural Networks*) until people managed to make it work (*Tim Cooijmans et al, 2016: Recurrent Batch Normalization*; specifically, initializing $\gamma = 0.1$ did the trick).



(a) We visualize the gradient flow through a batch-normalized tanh RNN as a function of $\gamma$. High variance causes vanishing gradient.

(b) We show the empirical expected derivative and interquartile range of tanh nonlinearity as a function of input variance. High variance causes saturation, which decreases the expected derivative.

*Figure 1 of "Recurrent Batch Normalization", https://arxiv.org/abs/1603.09025*

## Batch Normalization

Neuron value is normalized across the minibatch, and in case of CNN also across all positions.

## Layer Normalization

Neuron value is normalized across the layer.

*tolike se poruzlin' u RNN nyicisteji*

*kdyz nepmaiju pres batch, tahie se s tii blip pucyje*



Batch Norm          Layer Norm          Instance Norm          **Group Norm**

Figure 2 of "Group Normalization", https://arxiv.org/abs/1803.08494

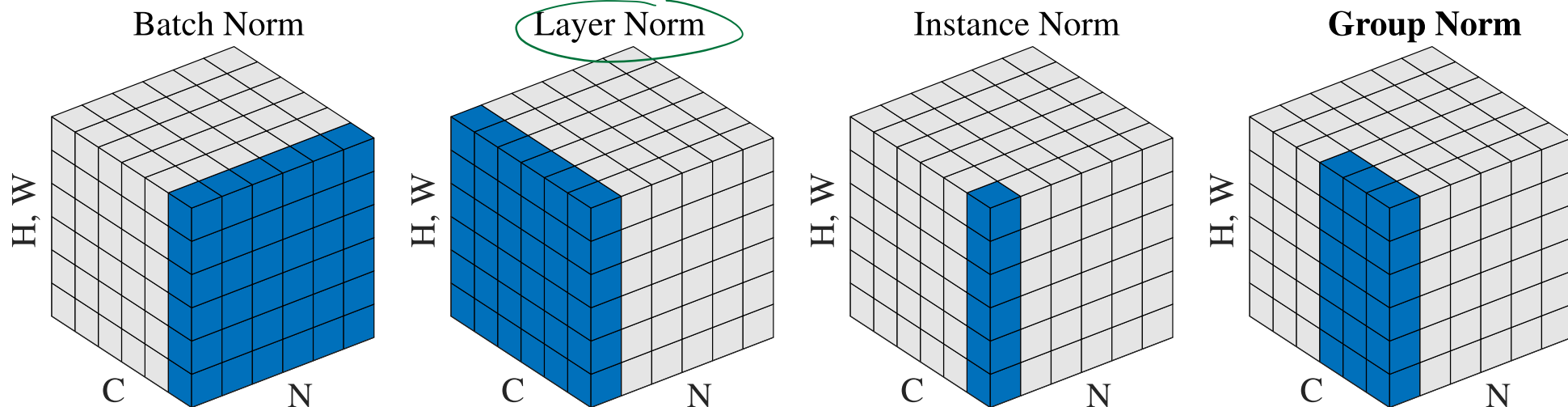# Layer Normalization

Consider a hidden value $\boldsymbol{x} \in \mathbb{R}^D$. Layer normalization (both during training and during inference) is performed as follows.

**Inputs**: An example $\boldsymbol{x} \in \mathbb{R}^D$, $\varepsilon \in \mathbb{R}$ with default value 0.001
**Parameters**: $\boldsymbol{\beta} \in \mathbb{R}^D$ initialized to $\mathbf{0}$, $\boldsymbol{\gamma} \in \mathbb{R}^D$ initialized to $\mathbf{1}$
**Outputs**: Normalized example $\boldsymbol{y}$

- $\mu \leftarrow \frac{1}{D} \sum_{i=1}^{D} x_i$
- $\sigma^2 \leftarrow \frac{1}{D} \sum_{i=1}^{D} (x_i - \mu)^2$
- $\hat{\boldsymbol{x}} \leftarrow (\boldsymbol{x} - \mu)/\sqrt{\sigma^2 + \varepsilon}$
- $\boldsymbol{y} \leftarrow \boldsymbol{\gamma} \odot \hat{\boldsymbol{x}} + \boldsymbol{\beta}$

## Layer Normalization

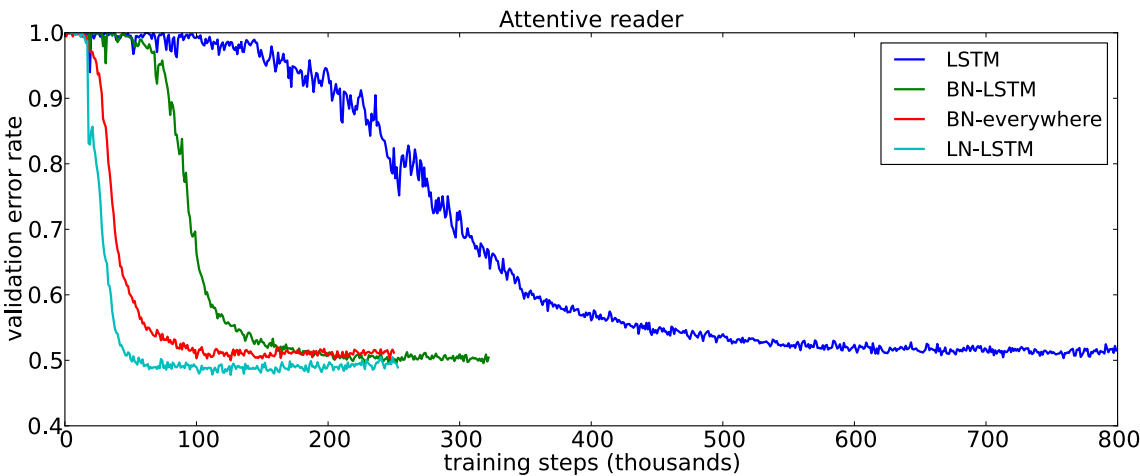Much more stable than batch normalization for RNN regularization.



Figure 2: Validation curves for the attentive reader model. BN results are taken from [Cooijmans et al., 2016].
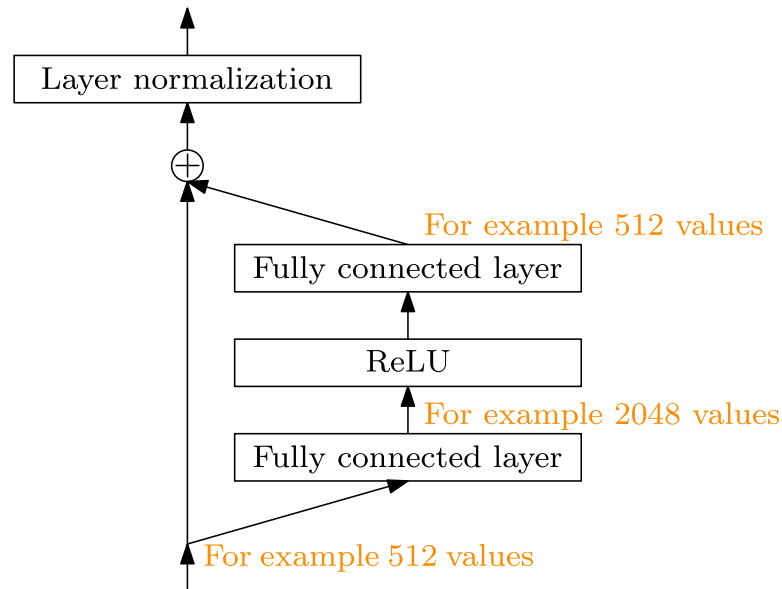
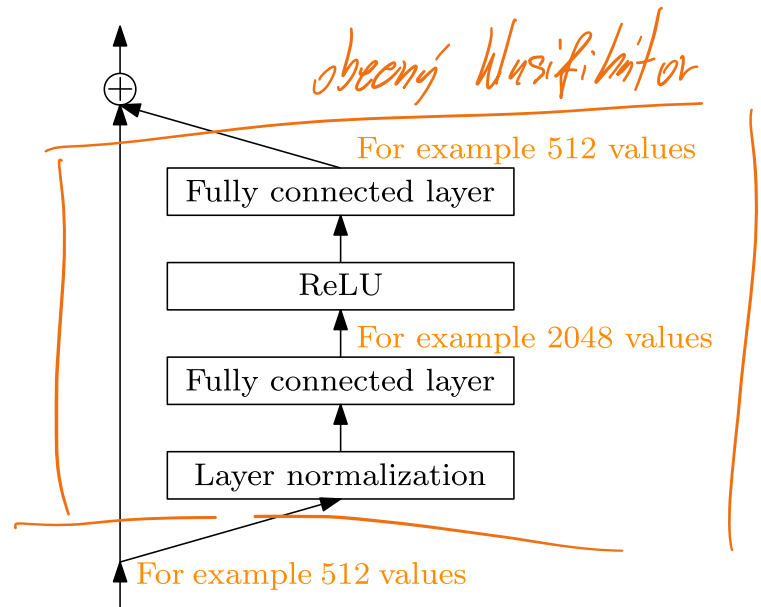| | Weight matrix re-scaling | Weight matrix re-centering | Weight vector re-scaling | Dataset re-scaling | Dataset re-centering | Single training case re-scaling |
|---|---|---|---|---|---|---|
| Batch norm | Invariant | No | Invariant | Invariant | Invariant | No |
| Weight norm | Invariant | No | Invariant | No | No | No |
| Layer norm | Invariant | Invariant | No | Invariant | No | Invariant |

# Layer Normalization

In an important recent architecture (namely Transformer), many fully connected layers are used, with a residual connection and a layer normalization.

Original "Post-LN" configuration
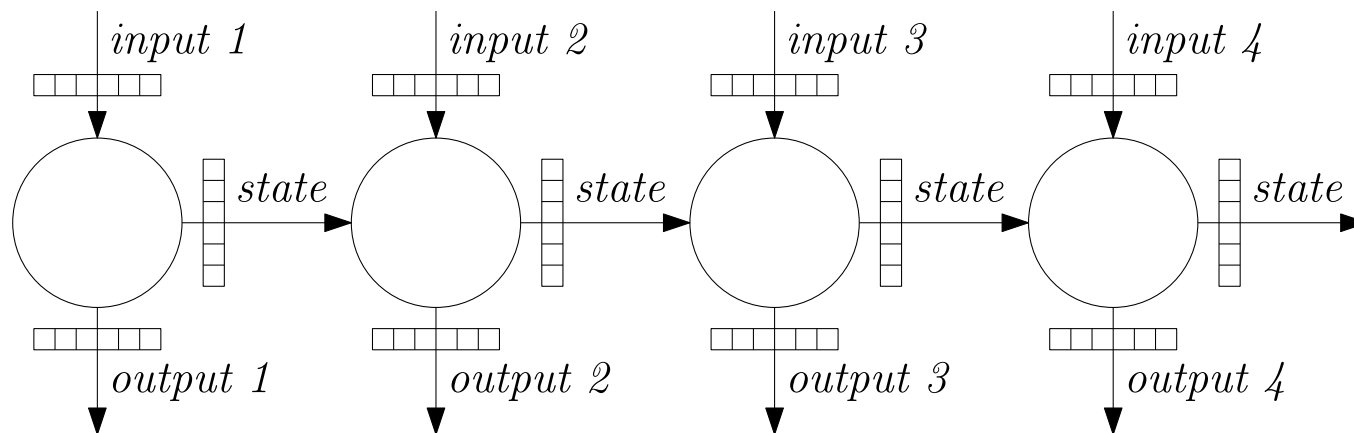
Improved "Pre-LN" configuration since 2020



This could be considered an alternative to highway networks, i.e., a suitable residual connection for fully connected layers. Note the architecture can be considered as a variant of a mobile inverted bottleneck $1 \times 1$ convolution block.

## Sequence Element Representation    *Contextualizovaná reprezentace*

Create output for individual elements, for example for classification of the individual elements.
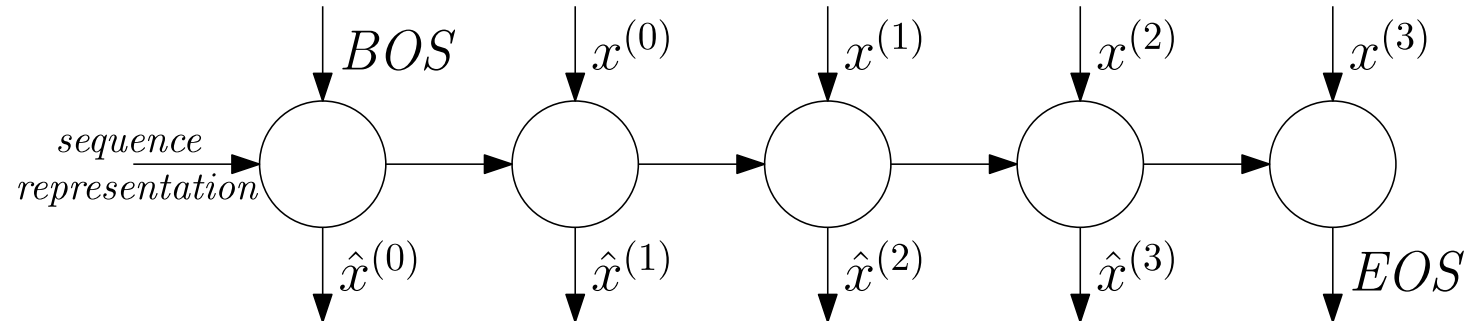


## Sequence Representation

Generate a single output for the whole sequence (either the last output or the last state).

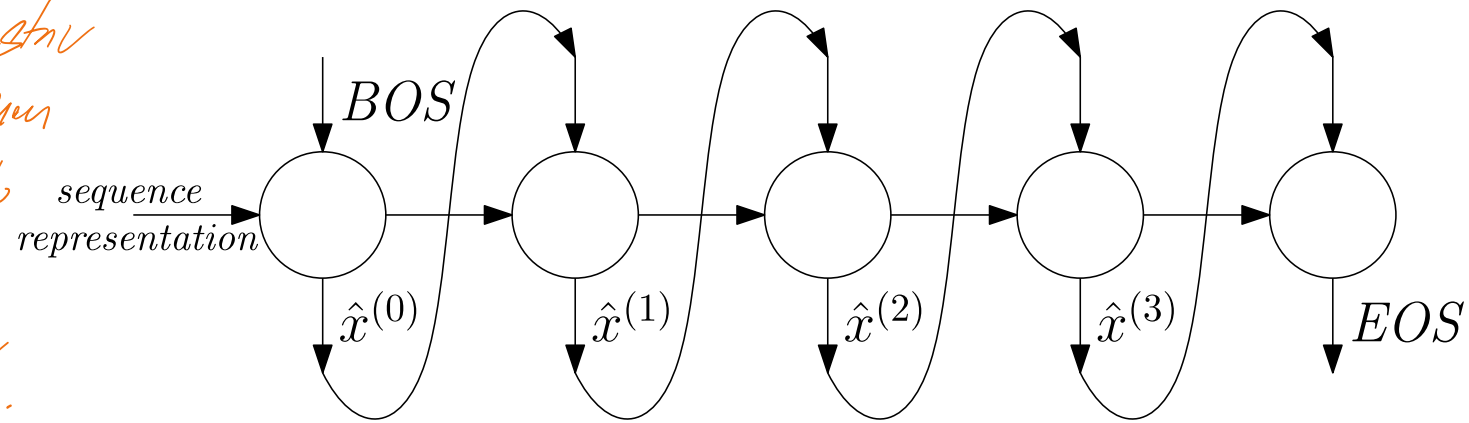*— pokud ale chci, aby reprezentace brala v potaz i historii*

## Sequence Prediction

During training, predict next sequence element.



During inference, use predicted elements as further inputs.

*sice sko), tý stav*
*předcho podobrou*
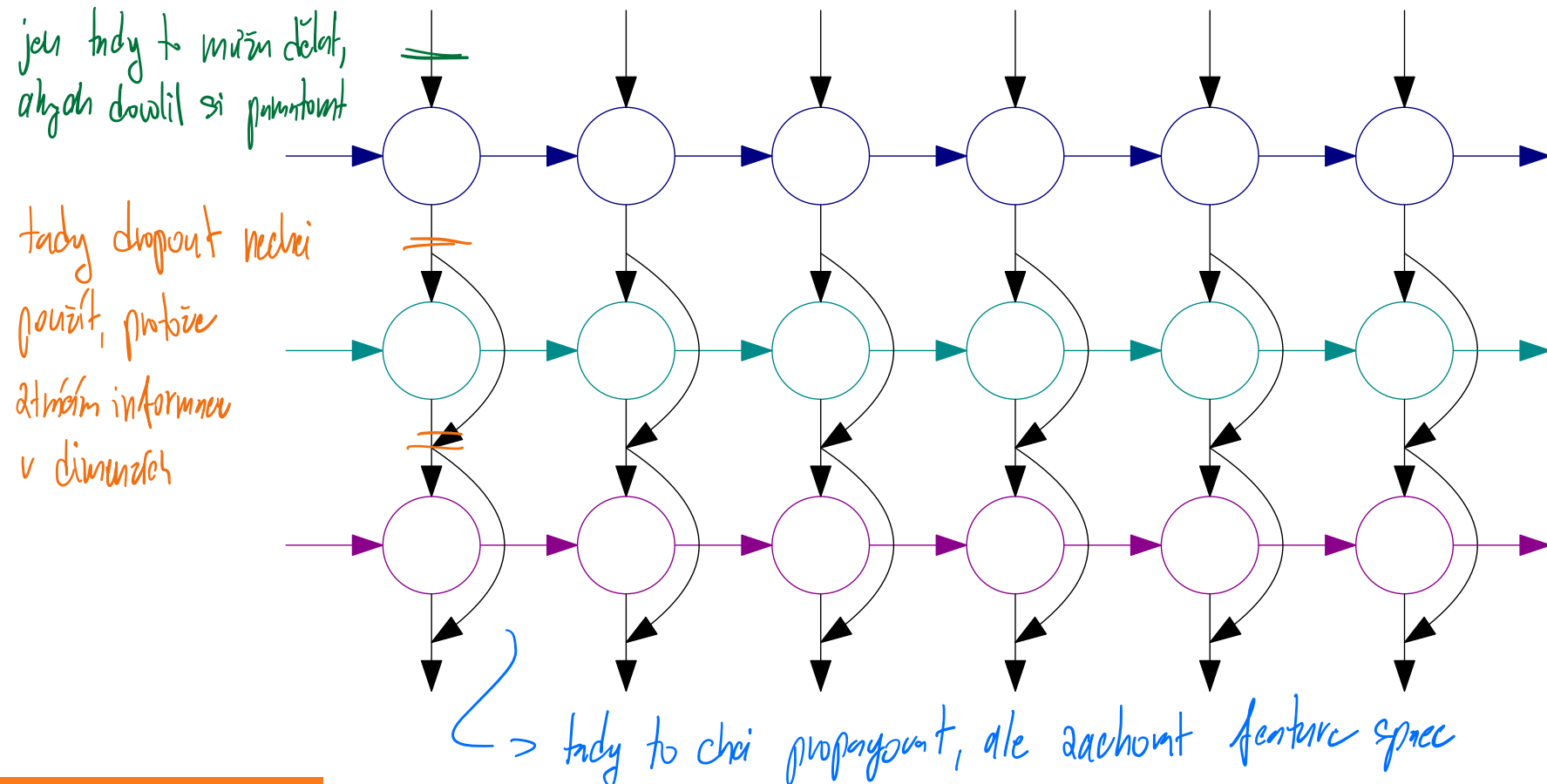*informaci, přesto*
*se tj vstupy*
*fakule přepojují.*

# Multilayer RNNs

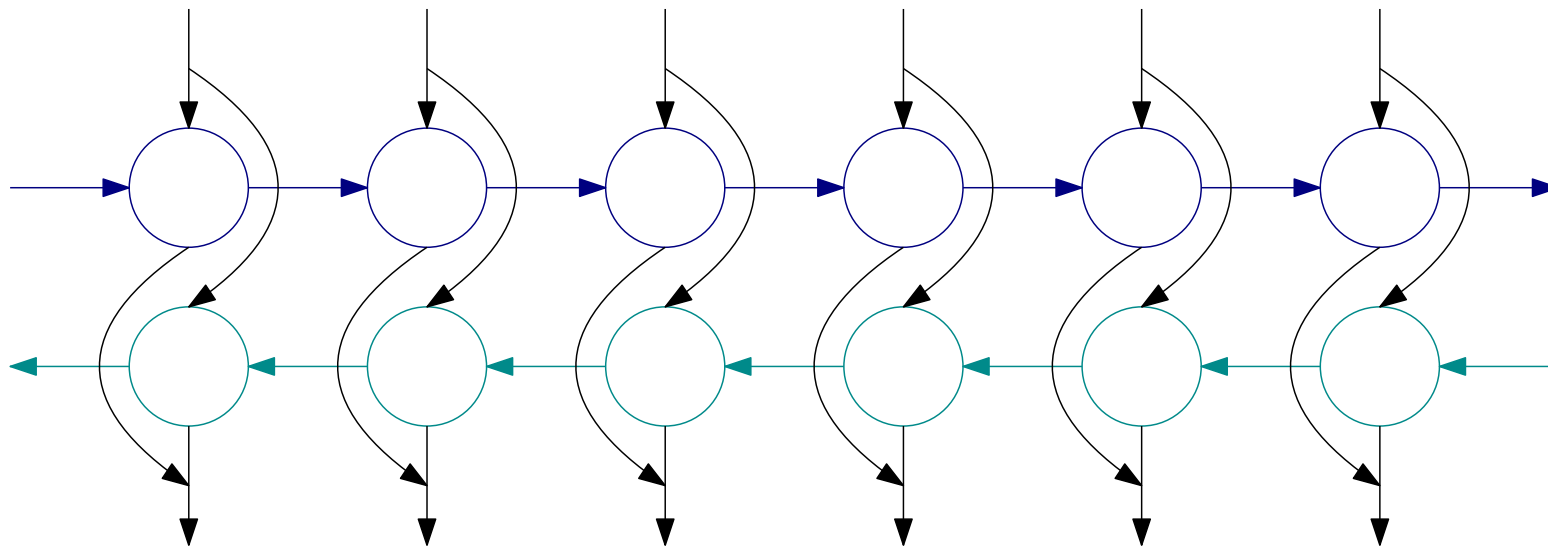We might stack several layers of recurrent neural networks. Usually using two or three layers gives better results than just one.

In case of multiple layers, residual connections usually improve results. Because dimensionality has to be the same, they are usually applied from the second layer.

To consider both the left and right contexts, a **bidirectional** RNN can be used, which consists of parallel application of a **forward** RNN and a **backward** RNN.



The outputs of both directions can be either **added** or **concatenated**. Even if adding them does not seem very intuitive, it does not increase dimensionality and therefore allows residual connections to be used in case of multilayer bidirectional RNN.

# Word Embeddings

We might represent **words** using one-hot encoding, considering 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.

These embeddings are initialized randomly and trained together with the rest of the network.

— třeba „pes" bude reprezentovaný položkami jako: „domácí zvíře", „Snave" atd...

tedy nějakými obecnými abstraktními virtuálními vlastnostmi

# Word Embeddings

The word embedding layer is in fact just a fully connected layer on top of one-hot encoding. However, it is not implemented in that way.

Instead, the so-called **embedding** layer is used, which is much more efficient. When a matrix is multiplied by an one-hot encoded vector (all but one zeros and exactly one 1), the row corresponding to that 1 is selected, so the embedding layer can be implemented only as a simple lookup.

*kolik máme slov*

*→ jak je chci reprezentovat*

In Keras, the embedding layer is available as

```
keras.layers.Embedding(input_dim, output_dim)
```

In PyTorch, it is available as

```
torch.nn.Embedding(input_dim, output_dim)
```

*↳ LSTM potřebuje 4 matice, tohže 4 obří matice pro fully-connected vrstvy. Zatímco embeding mi vrátí jen nějaký lmitkf vektor.*

Even if the embedding layer is just a fully connected layer on top of one-hot encoding, it is important that this layer is *shared* across the whole network.

*Blbý ale je, že existují slova, co jsem při tréningu neviděl a neumím je embedingovat*

## Recurrent Character-level WEs

In order to handle words not seen during training, we could find a way to generate a representation from the word **characters**.

A possible way to compose the representation from individual characters is to use RNNs – we embed *characters* to get character representation, and then use an RNN to produce the representation of a whole *sequence of characters*.

Usually, both forward and backward directions are used, and the resulting representations are concatenated/added.

*Rekurentní síť!*


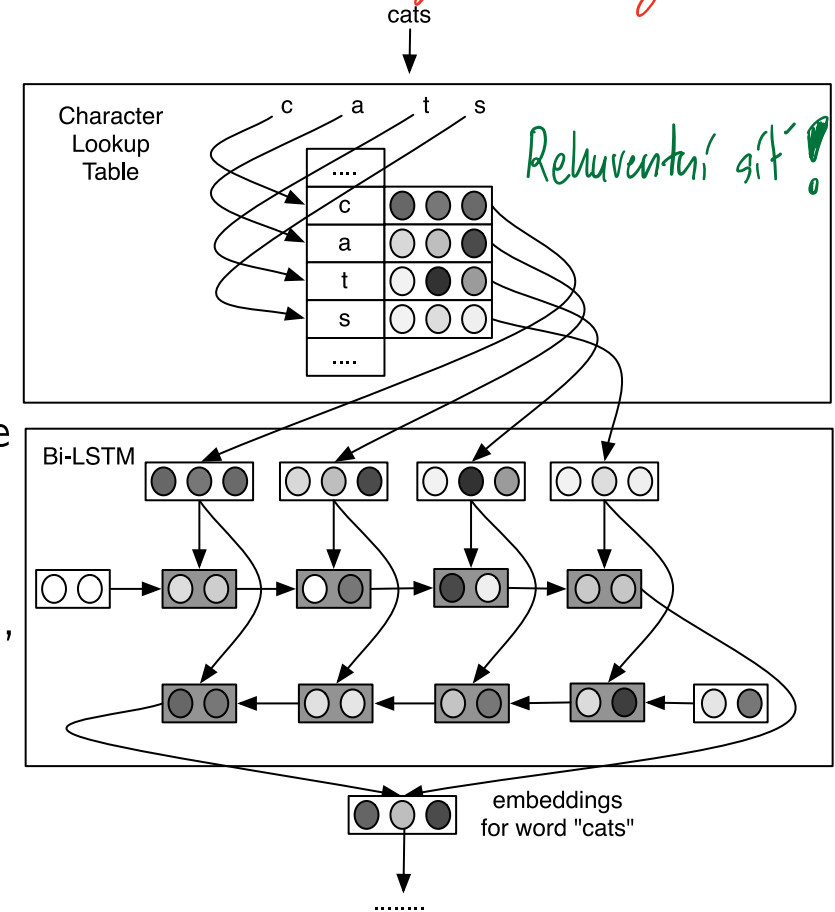
Figure 1 of "Finding Function in Form: Compositional Character Models for Open Vocabulary Word Representation", https://arxiv.org/abs/1508.02096

*Paměť buněk drží obecný význam slov. Lepší než dostat nulový vektor...*

## Convolutional Character-level WEs

Alternatively, 1D convolutions might be used.

Assume we use a 1D convolution with kernel size 3. It produces a representation for every input word trigram, but we need a representation of the whole word. To that end, we use *global max-pooling* — using it has an interpretable meaning, where the kernel is a *pattern* and the activation after the maximum is a level of a highest match of the pattern anywhere in the word.

Kernels of varying sizes are usually used (because it makes sense to have patterns for unigrams, bigrams, trigrams, ...) — for example, 25 filters for every kernel size $(1, 2, 3, 4, 5)$ might be used.

Lastly, authors employed a highway layer after the convolutions, improving the results (compared to not using any layer or using a fully connected one).
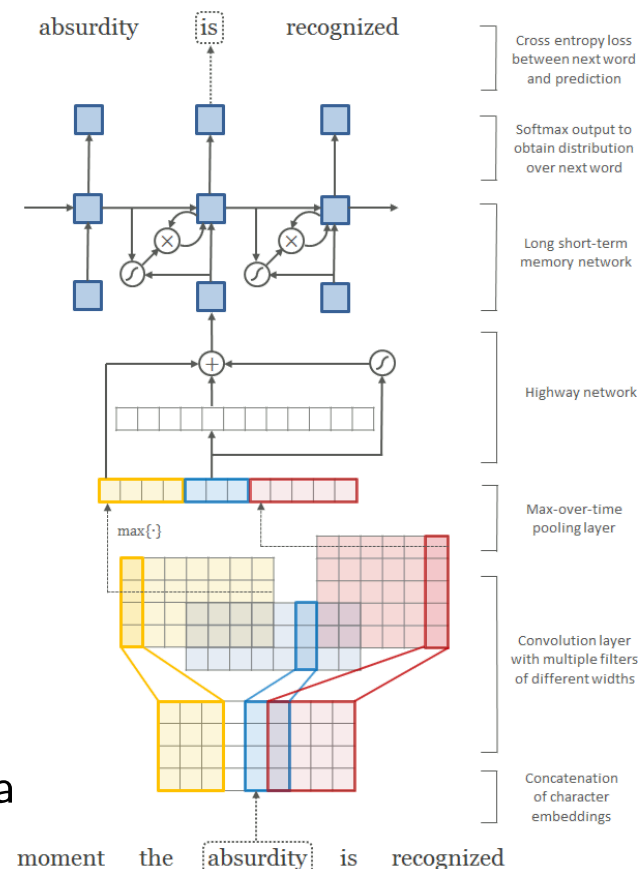


*Figure 1 of "Character-Aware Neural Language Models",*
*https://arxiv.org/abs/1508.06615*

| *increased* | *John* | *Noahshire* | *phding* |
|---|---|---|---|
| reduced | Richard | Nottinghamshire | mixing |
| improved | George | Bucharest | modelling |
| expected | James | Saxony | styling |
| decreased | Robert | Johannesburg | blaming |
| targeted | Edward | Gloucestershire | christening |

Table 2: Most-similar in-vocabular words under the C2W model; the two query words on the left are in the training vocabulary, those on the right are nonce (invented) words.

Table 2 of "Finding Function in Form: Compositional Character Models for Open Vocabulary Word Representation", https://arxiv.org/abs/1508.02096

| | In Vocabulary | | | | | Out-of-Vocabulary | | |
|---|---|---|---|---|---|---|---|---|
| | *while* | *his* | *you* | *richard* | *trading* | *computer-aided* | *misinformed* | *looooook* |
| LSTM-Word | although | your | conservatives | jonathan | advertised | – | – | – |
| | letting | her | we | robert | advertising | – | – | – |
| | though | my | guys | neil | turnover | – | – | – |
| | minute | their | i | nancy | turnover | – | – | – |
| LSTM-Char (before highway) | chile | this | your | hard | heading | computer-guided | informed | look |
| | whole | hhs | young | rich | training | computerized | performed | cook |
| | meanwhile | is | four | richer | reading | disk-drive | transformed | looks |
| | white | has | youth | richter | leading | computer | inform | shook |
| LSTM-Char (after highway) | meanwhile | hhs | we | eduard | trade | computer-guided | informed | look |
| | whole | this | your | gerard | training | computer-driven | performed | looks |
| | though | their | doug | edward | traded | computerized | outperformed | looked |
| | nevertheless | your | i | carl | trader | computer | transformed | looking |

**Table 6:** Nearest neighbor words (based on cosine similarity) of word representations from the large word-level and character-level (before and after highway layers) models trained on the PTB. Last three words are OOV words, and therefore they do not have representations in the word-level model.

*Table 6 of "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

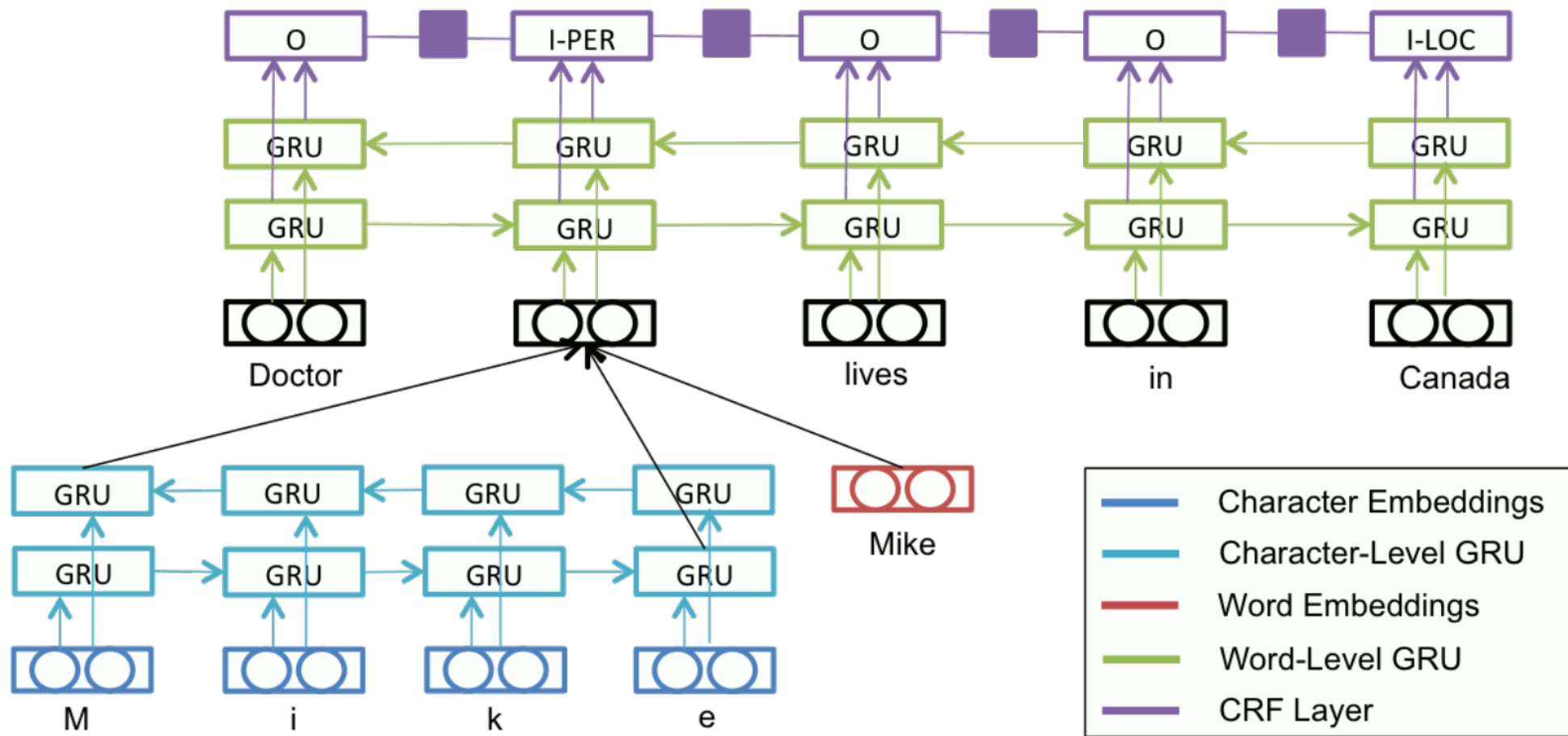- We can cache character-level word embeddings during inference.

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