#### NPFL114, Lecture 08

# Recurrent Neural Networks II, Word Embeddings



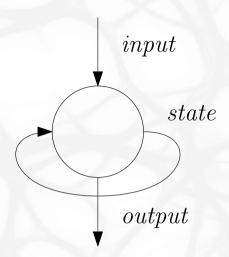


Milan Straka

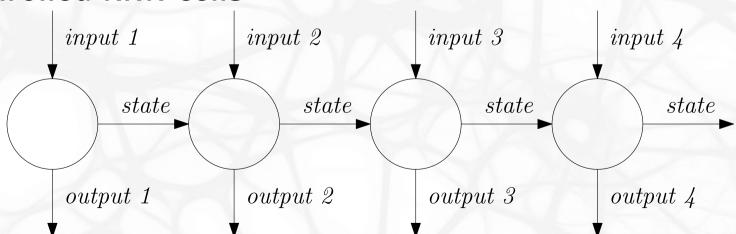
#### **Recurrent Neural Networks**



#### Single RNN cell



#### **Unrolled RNN cells**



## Challenge of Long-term Dependencies



Consider a RNN cell which given an input  ${\pmb x}^{(t)}$  and previous state  ${\pmb s}^{(t-1)}$  computes the new state as

$$oldsymbol{s}^{(t)} = f(oldsymbol{s}^{(t-1)}, oldsymbol{x}^{(t)}; oldsymbol{ heta}).$$

## Challenge of Long-term Dependencies



Consider a RNN cell which given an input  $m{x}^{(t)}$  and previous state  $m{s}^{(t-1)}$  computes the new state as

$$oldsymbol{s}^{(t)} = f(oldsymbol{s}^{(t-1)}, oldsymbol{x}^{(t)}; oldsymbol{ heta}).$$

RNN cells generally suffer a lot from vanishing/exploding gradients (the challenge of long-term dependencies), a problem attributed to the fact that same function is iteratively applied many times.

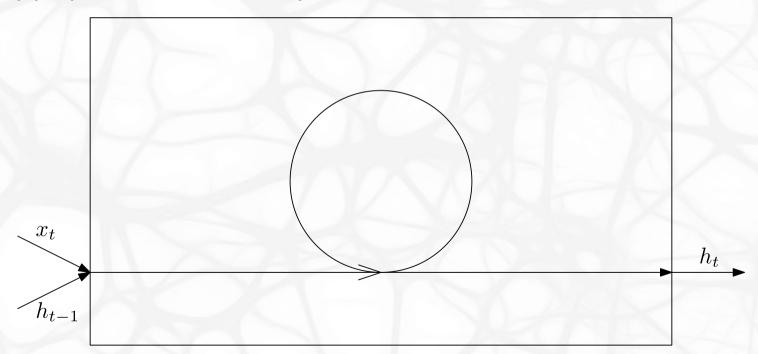
## Long Short-Term Memory



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

$$f'=1$$
.

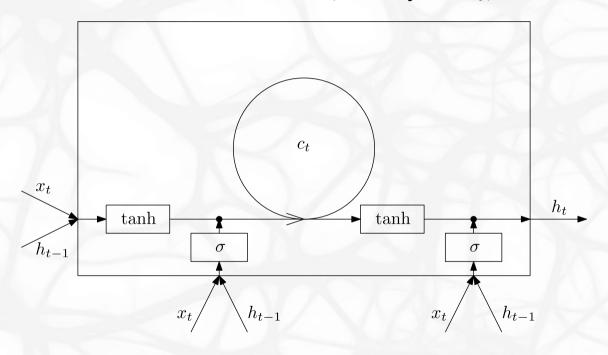
They propose to achieve that by a constant error carrousel.



#### 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*  $c_t$ ).

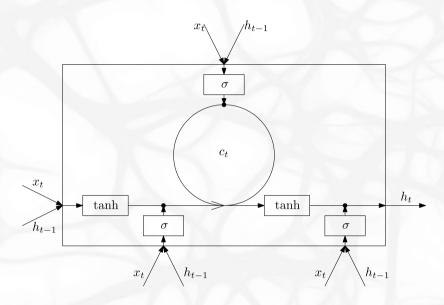


$$egin{aligned} m{i}_t &\leftarrow \sigma(m{W}^im{x}_t + m{V}^im{h}_{t-1} + m{b}^i) \ m{o}_t &\leftarrow \sigma(m{W}^om{x}_t + m{V}^om{h}_{t-1} + m{b}^o) \ m{c}_t &\leftarrow m{c}_{t-1} + m{i}_t \cdot anh(m{W}^ym{x}_t + m{V}^ym{h}_{t-1} + m{b}^y) \ m{h}_t &\leftarrow m{o}_t \cdot anh(m{c}_t) \end{aligned}$$

#### Long Short-Term Memory



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



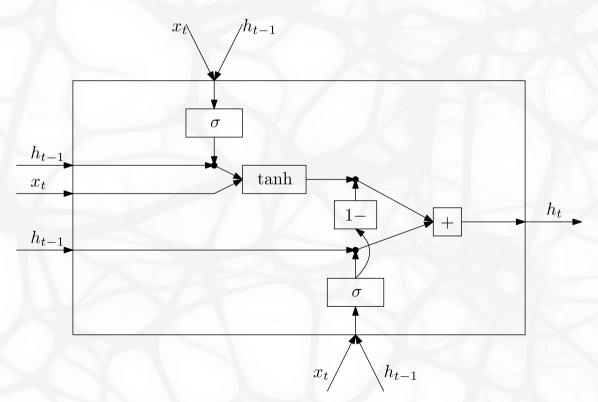
$$egin{aligned} oldsymbol{i}_t &\leftarrow \sigma(oldsymbol{W}^i oldsymbol{x}_t + oldsymbol{V}^i oldsymbol{h}_{t-1} + oldsymbol{b}^i) \ oldsymbol{f}_t &\leftarrow \sigma(oldsymbol{W}^f oldsymbol{x}_t + oldsymbol{V}^f oldsymbol{h}_{t-1} + oldsymbol{b}^o) \ oldsymbol{c}_t &\leftarrow oldsymbol{f}_t \cdot oldsymbol{c}_{t-1} + oldsymbol{i}_t \cdot anh(oldsymbol{W}^y oldsymbol{x}_t + oldsymbol{V}^y oldsymbol{h}_{t-1} + oldsymbol{b}^y) \ oldsymbol{h}_t &\leftarrow oldsymbol{o}_t \cdot anh(oldsymbol{c}_t) \end{aligned}$$

#### **Gated Recurrent Unit**



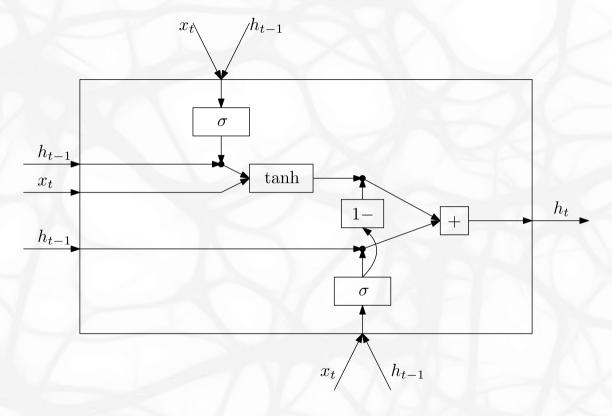
Gated recurrent unit (GRU) was proposed by Cho et al. (2014) as a simplification of LSTM. The main differences are

- no memory cell
- forgetting and updating tied together



#### **Gated Recurrent Unit**





$$egin{aligned} oldsymbol{r}_t &\leftarrow \sigma(oldsymbol{W}^roldsymbol{x}_t + oldsymbol{V}^roldsymbol{h}_{t-1} + oldsymbol{b}^r) \ oldsymbol{u}_t &\leftarrow \sigma(oldsymbol{W}^uoldsymbol{x}_t + oldsymbol{V}^uoldsymbol{h}_{t-1} + oldsymbol{b}^u) \ oldsymbol{\hat{h}}_t &\leftarrow anh(oldsymbol{W}^holdsymbol{x}_t + oldsymbol{r}_t \cdot oldsymbol{V}^holdsymbol{h}_{t-1} + oldsymbol{b}^u) \ oldsymbol{h}_t &\leftarrow oldsymbol{u}_t \cdot oldsymbol{h}_{t-1} + (1 - oldsymbol{u}_t) \cdot oldsymbol{\hat{h}}_t \end{aligned}$$



**Word Embeddings** 



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.



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.



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.

## **Unsupervised Word Embeddings**



The embeddings can be trained for each task separately.

#### **Unsupervised Word Embeddings**



The embeddings can be trained for each task separately.

However, a method of precomputing word embeddings have been proposed, based on *distributional hypothesis*:

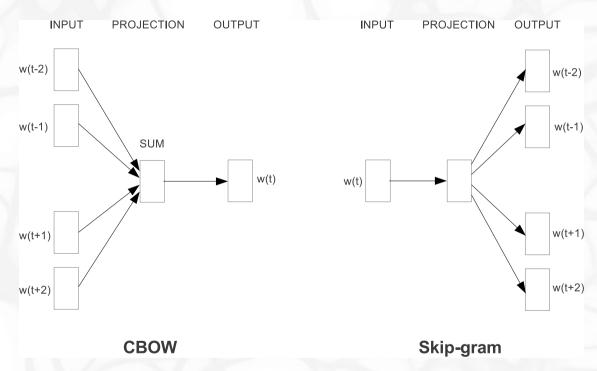
Words that are used in the same contexts tend to have similar meanings.

The distributional hypothesis is usually attributed to Firth (1957).

#### Word2Vec



Mikolov et al. (2013) proposed two very simple architectures for precomputing word embeddings, together with a C multi-threaded implementation word2vec.



#### Word2Vec

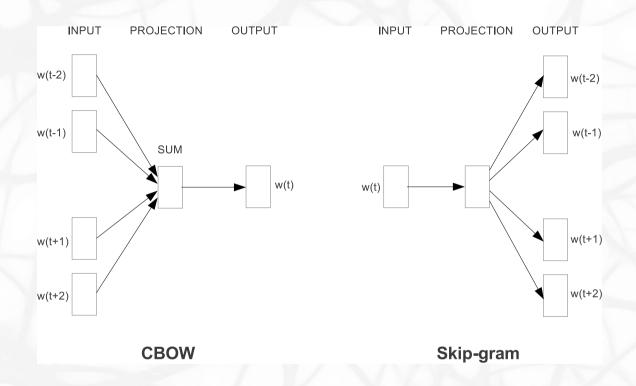


Table 8: Examples of the word pair relationships, using the best word vectors from Table 4 (Skipgram model trained on 783M words with 300 dimensionality).

Relationship	Relationship Example 1		Example 3	
France - Paris Italy: Rome		Japan: Tokyo	Florida: Tallahassee	
big - bigger	g - bigger small: larger		quick: quicker	
Miami - Florida	Miami - Florida Baltimore: Maryland		Kona: Hawaii	
Einstein - scientist Messi: midfield		Mozart: violinist	Picasso: painter	
Sarkozy - France	Berlusconi: Italy	Merkel: Germany	Koizumi: Japan	
copper - Cu	zinc: Zn	gold: Au	uranium: plutonium	
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack	
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone	
Microsoft - Ballmer	Microsoft - Ballmer Google: Yahoo		Apple: Jobs	
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza	

#### Word2Vec - SkipGram Model





Considering input word  $w_i$  and output  $w_o$ , the Skip-gram model defines

$$p(w_o|w_i) \stackrel{ ext{def}}{=} rac{e^{ extbf{ extit{W}}_{w_o}^ op extbf{ extit{V}}_{w_i}}}{\sum_w e^{ extbf{ extit{W}}_w^ op extbf{ extit{V}}_{w_i}}}.$$

#### Word2Vec - Hierarchical Softmax



Instead of a large softmax, we construct a binary tree over the words, with a sigmoid classifier for each node.

If word w corresponds to a path  $n_1, n_2, \ldots, n_L$ , we define

$$p_{ ext{HS}}(w|w_i) \stackrel{ ext{def}}{=} \prod_{j=1}^{L-1} \sigma([+1 ext{ if } n_{j+1} ext{ is right child else -1}] \cdot oldsymbol{W}_{n_j}^ op oldsymbol{V}_{w_i}).$$

### Word2Vec - Negative Sampling



Instead of a large softmax, we could train individual sigmoids for all words.

We could also only sample the *negative examples* instead of training all of them.

This gives rise to the following *negative sampling* objective:

$$l_{ ext{NEG}}(w_o, w_i) \stackrel{ ext{ iny def}}{=} \log \sigma(oldsymbol{W}_{w_o}^ op oldsymbol{V}_{w_i}) - \sum_{j=1}^k \mathbb{E}_{w_j \sim P(w)} \log \sigma(-oldsymbol{W}_{w_j}^ op oldsymbol{V}_{w_i}).$$

### Word2Vec - Negative Sampling



Instead of a large softmax, we could train individual sigmoids for all words.

We could also only sample the *negative examples* instead of training all of them.

This gives rise to the following *negative sampling* objective:

$$l_{ ext{NEG}}(w_o, w_i) \stackrel{ ext{ iny def}}{=} \log \sigma(oldsymbol{W}_{w_o}^ op oldsymbol{V}_{w_i}) - \sum_{j=1}^k \mathbb{E}_{w_j \sim P(w)} \log \sigma(-oldsymbol{W}_{w_j}^ op oldsymbol{V}_{w_i}).$$

For P(w), both uniform and unigram distribution U(w) work, but

$$U(w)^{3/4}$$

outperforms them significantly.

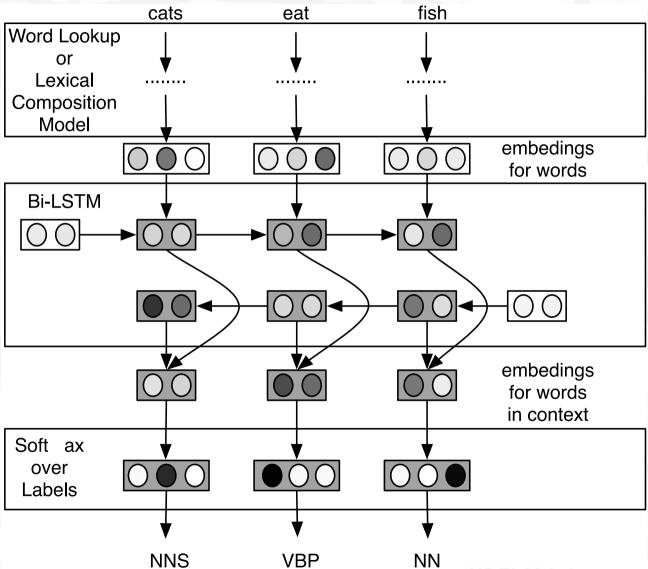
## **Basic NLP Processing**



**Basic NLP Processing** 

#### **Bidirectional RNN**





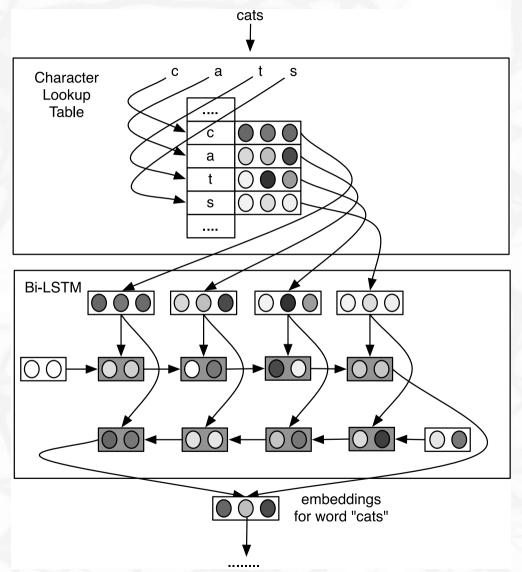
## Word Embeddings for Unknown Words



Word Embeddings for Unknown Words

#### Recurrent Character-level WEs





#### Recurrent Character-level WEs

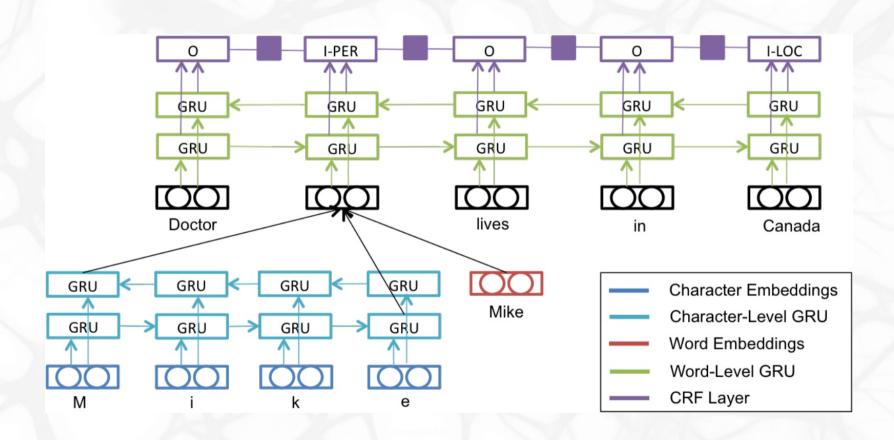


increased	John	Noahshire	phding	
reduced	Richard	Nottinghamshire	mixing	
improved	George	Bucharest	modelling	
expected	James	Saxony	styling	
decreased	Robert	Johannesburg	blaming	
targeted	Edward	Gloucestershire	christening	

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.

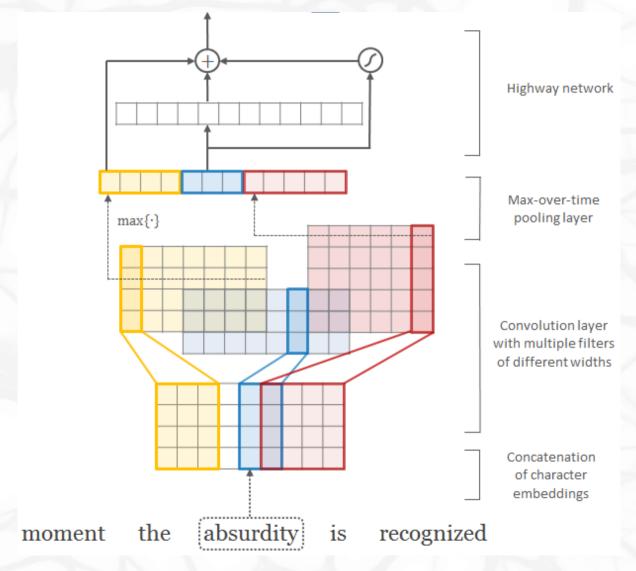
#### Recurrent Character-level WEs





#### Convolutional Character-level WEs





#### Convolutional Character-level WEs



In Vocabulary				Out-of-Vocabulary			
while	his	you	richard	trading	computer-aided	misinformed	looooook
although	your	conservatives	jonathan	advertised			_
letting	her	we	robert	advertising	-	_	_
though	my	guys	neil	turnover	_	_	-
minute	their	i	nancy	turnover		-	_
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
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
	although letting though minute  chile whole meanwhile white  meanwhile whole though	although your letting her though my minute their  chile this whole hhs meanwhile is white has  meanwhile hhs whole this though their	while his you  although your conservatives letting her we though my guys minute their i  chile this your whole hhs young meanwhile is four white has youth  meanwhile hhs we whole this your though their doug	while his you richard  although your conservatives jonathan letting her we robert though my guys neil minute their i nancy  chile this your hard whole hhs young rich meanwhile is four richer white has youth richter  meanwhile hhs we eduard whole this your gerard though their doug edward	whilehisyourichardtradingalthoughyourconservativesjonathanadvertisedlettingherwerobertadvertisingthoughmyguysneilturnoverminutetheirinancyturnoverchilethisyourhardheadingwholehhsyoungrichtrainingmeanwhileisfourricherreadingwhitehasyouthrichterleadingmeanwhilehhsweeduardtradewholethisyourgerardtrainingthoughtheirdougedwardtraded	while his you richard trading computer-aided  although your conservatives jonathan advertised — letting her we robert advertising — though my guys neil turnover — minute their i nancy turnover —  chile this your hard heading computer-guided whole hhs young rich training computerized meanwhile is four richer reading disk-drive white has youth richter leading computer  meanwhile hhs we eduard trade computer-guided whole this your gerard training computer-driven though their doug edward traded computerized	while his you richard trading computer-aided misinformed  although your conservatives jonathan advertised — — —  letting her we robert advertising — — —  though my guys neil turnover — — —  minute their i nancy turnover — — —  chile this your hard heading computer-guided informed whole hhs young rich training computerized performed meanwhile is four richer reading disk-drive transformed white has youth richter leading computer inform  meanwhile hhs we eduard trade computer-guided informed training computer inform

**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.

### Character N-grams



Another simple idea appeared simultaneously in three nearly simultaneous publications as <a href="Charagram">Charagram</a>, <a href="Subword Information">SubGram</a>.

A word embedding is a sum of the word embedding plus embeddings of its character *n*-grams. Such embedding can be pretrained using same algorithms as word2vec.

### Character N-grams



Another simple idea appeared simultaneously in three nearly simultaneous publications as <a href="Charagram">Charagram</a>, <a href="Subword Information">SubGram</a>.

A word embedding is a sum of the word embedding plus embeddings of its character *n*-grams. Such embedding can be pretrained using same algorithms as word2vec.

The implementation can be

- dictionary based: only some number of frequent character n-grams is kept
- ullet hash-based: character n-grams are hashed into K buckets (usually  $K\sim 10^6$  is used)

## **Charagram WEs**



query	tiling	tech-rich	english-born	micromanaging	eateries	dendritic
sisg	tile flooring	tech-dominated tech-heavy	british-born polish-born	micromanage micromanaged	restaurants eaterie	dendrites
sg	bookcases built-ins	technology-heavy .ixic	most-capped ex-scotland	defang internalise	restaurants delis	epithelial p53

Table 7: Nearest neighbors of rare words using our representations and skipgram. These hand picked examples are for illustration.

### **Charagram WEs**



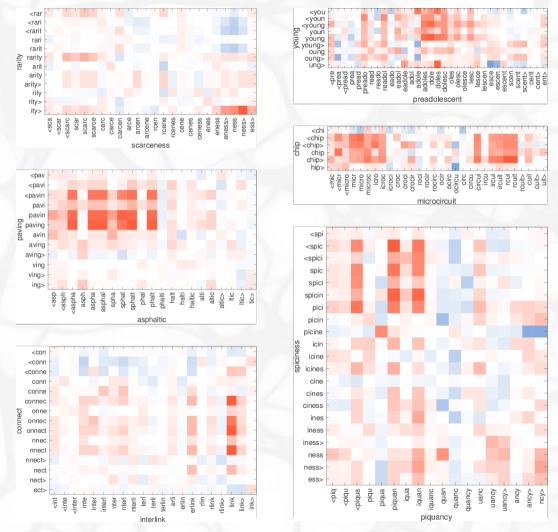


Figure 2: Illustration of the similarity between character n-grams in out-of-vocabulary words. For each pair, only one word is OOV, and is shown on the x axis. Red indicates positive cosine, while blue negative.

# Character-level WE Implementation



#### **Training**

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

## Character-level WE Implementation



#### **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.

# TF – RNN and Variable Length Sequences



#### tf.nn.dynamic\_rnn

## TF - RNN and Variable Length Sequences



#### tf.nn.dynamic\_rnn

The sequence\_length parameter allows sequences of different length. If used, the outputs past the sequence lengths are zeros and the state is the state after processing last element of each sequence.

## TF - RNN and Variable Length Sequences



#### Losses

```
outputs, state = tf.nn.sparse_softmax_cross_entropy(
   labels,
   logits,
   weights=1.0,
   ...)
```

## TF - RNN and Variable Length Sequences



#### Losses

```
outputs, state = tf.nn.sparse_softmax_cross_entropy(
   labels,
   logits,
   weights=1.0,
   ...)
```

The weights can be used to mask some of the labels and logits.

#### Generating the weights

```
tf.sequence_masks(lengths, maxlen=None, dtype=tf.bool)
weights = tf.sequence_masks(lengths, dtype=tf.float32)
```

#### TF - Bidirectional RNNs



#### tf.nn.bidirectional\_dynamic\_rnn





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

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



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

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

Highway networks add residual connection with gating:

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



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

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

Highway networks add residual connection with gating:

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

Usually, the gating is defined as

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



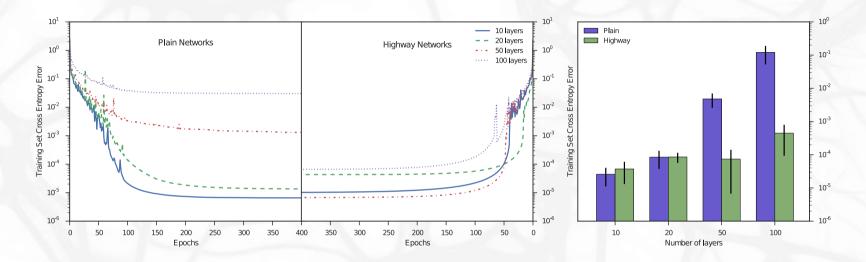
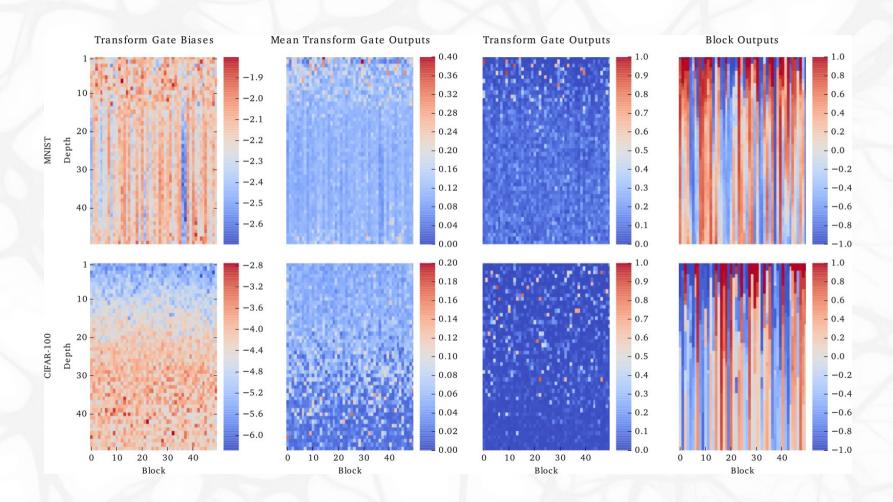


Figure 1: Comparison of optimization of plain networks and highway networks of various depths. *Left:* The training curves for the best hyperparameter settings obtained for each network depth. *Right:* Mean performance of top 10 (out of 100) hyperparameter settings. Plain n tworks become much harder 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).







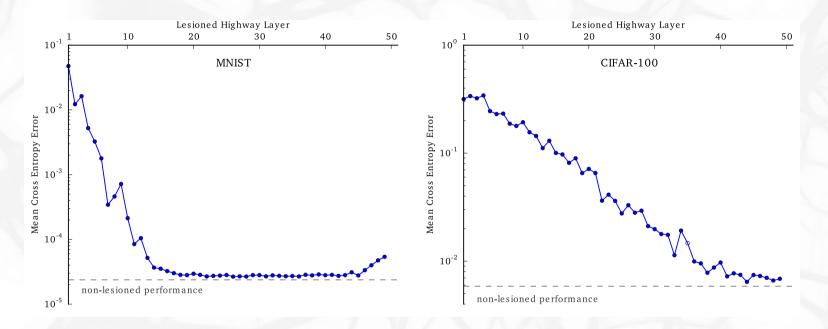


Figure 4: Lesioned training set performance (y-axis) of he best 50-layer highway networks on MNIST (left) and CIFAR-100 (right), as a function of the lesioned layer (x-axis). Evaluated on the full training set while forcefully closing all the transform gates of a single layer at a time. The non-lesioned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