

Natural Language Processing with Recurrent Neural Networks

2025-26

Ion Androutsopoulos

http://www.aueb.gr/users/ion/

Contents

- Recurrent neural networks (RNNs), GRUs/LSTMs.
- Bidirectional and stacked RNNs.
- RNNs with self-attention or global max pooling.
- RNNs in text and token classification, RNN language models.
- Obtaining word embeddings from character-based RNNs.
- Hierarchical RNNs.
- Sequence-to-sequence RNN models with attention, and applications in machine translation.
- Variational dropout.
- Universal sentence encoders, LASER.
- Pretraining RNN language models, ELMo.

Extracting contract elements

THIS AGREEMENT is made the 15th day of October 2009 (The "Effective Date") BETWEEN:

- (1) Sugar 13 Inc., a corporation whose office is at James House, 42-50 Bond Street, London, EW2H TL ("Sugar");
- (2) **E2 UK Limited**, a limited company whose registered office is at 260 Bathurst Road, Yorkshire, SL3 4SA ("**Provider**").

RECITALS:

A. The Parties wish to enter into a framework agreement which will enable Sugar, from time to time, to [...]

B. [...]

NO THEREFORE IT IS AGREED AS FOLLOWS:

ARTICLE I - DEFINITIONS

"Sugar" shall mean: Sugar 13 Inc.

"Provider" shall mean: E2 UK Limited

"1933 Act" shall mean: Securities Act of 1933

ARTICLE II - TERMINATION

The Service Period will be for five (5) years from the Effective Date (The "Initial Term"). The agreement is considered to be terminated in October 16, 2014.

ARTICLE III - PAYMENT - FEES

During the service period monthly payments should occur. The estimated fees for the Initial Term are £154,800.

ARTICLE IV - GOVERNING LAW

This agreement shall be governed and construed in accordance with the Laws of England & Wales. Each party hereby irrevocably submits to the exclusive jurisdiction of the courts sitting in Northern London.

IN WITNESS WHEREOF, the parties have caused their respective duly authorized officers to execute this Agreement.

BY: George Fake Authorized Officer Sugar 13 Inc.

BY: Olivier Giroux CEO E2 UK LIMITED

Identify start/end dates, duration, contractors, amount, legislations refs, jurisdiction etc. Similar to Named Entity Recognition (NER).

I. Chalkidis, I. Androutsopoulos, A. Michos, "Extracting Contract Elements", ICAIL 2017, http://nlp.cs.aueb.gr/pubs/icail2017.pdf.

RNN-based Named Entity Recognizer

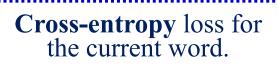
CrEnt

 $W^{(o)}$

 $W^{(e)}$

 $E\vec{x}_i$

 $W^{(h)}$



We can think of $W^{(o)}$ as containing class embeddings.

State: vector acting as a memory (remembering things about the words of the input seen so far).

1-hot vector of the current word.

Correct class probabilities for the current word ($|C| \times 1$, 1-hot).

Predicted probability distribution ($|C| \times 1$) over the classes (e.g., B-Person, I-Person, B-Location, Other) for the current word.

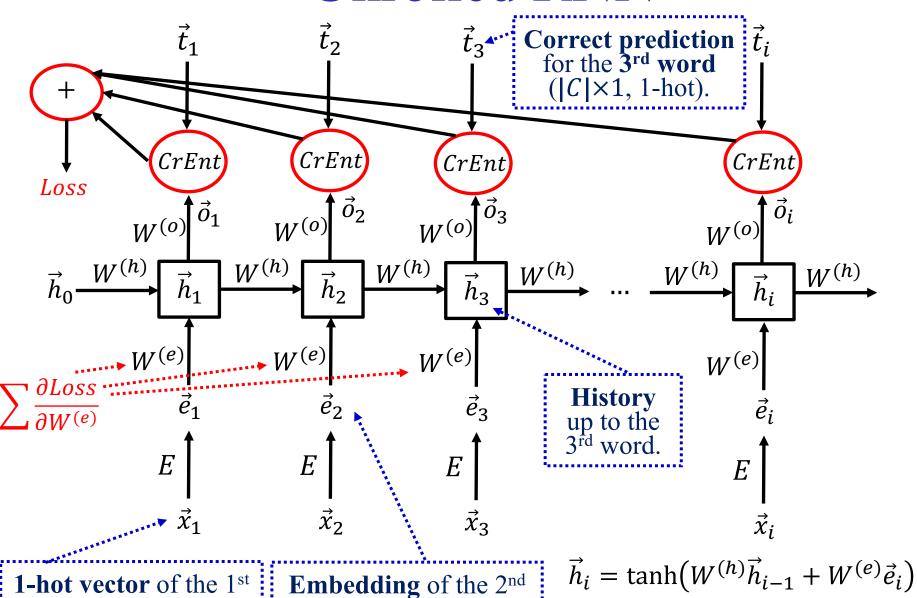
$$\vec{o}_i = \operatorname{softmax}(W^{(o)}\vec{h}_i)$$

 $\vec{h}_i = \tanh(W^{(h)}\vec{h}_{i-1} + W^{(e)}\vec{e}_i)$

The new state (memory) is a combination of the previous one and the new word embedding.

Embedding of the current word.

Unrolled RNN



word of the sentence

word of the sentence

 $\vec{o}_i = \operatorname{softmax}(W^{(o)}\vec{h}_i)$

4

RNN language model

 t_{i+1}

CrEnt

 $W^{(h)}$

 $W^{(e)}$

Cross-entropy loss for the prediction of the next word.

 $W^{(o)}$ contains **alternative** (output) word embeddings.

Some RNN LMs use E^T as $W^{(o)}$.

State: vector acting as a memory (remembering things about the words of the input seen so far).

1-hot vector of the current word.

Correct prediction for **next** word ($|V| \times 1$, 1-hot).

 $\vec{o}_{i+1} = \operatorname{softmax}(W^{(o)}\vec{h}_i)$ $W^{(o)}$

 $E\vec{x}_i$

Probability distribution ($|V| \times 1$) over the **vocabulary**. Shows which words the LM expects to see next.

 $\vec{h}_i = \tanh(W^{(h)}\vec{h}_{i-1} + W^{(e)}\vec{e}_i)$

The new state (memory) is a combination of the previous one and the new word embedding.

Embedding of the current word.

Reminder: LMs as next word predictors

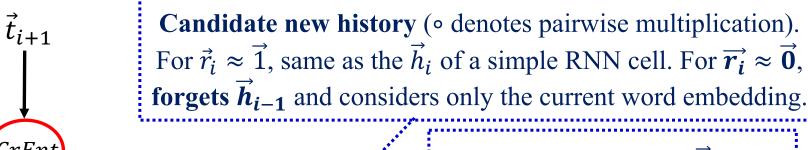
• Sequence probability using a bigram LM:

$$P(w_1^k) = P(w_1, ..., w_k) = P(w_1) \cdot P(w_2 \mid w_1) \cdot P(w_3 \mid w_1, w_2) \cdot P(w_4 \mid w_1^3) \cdots P(w_k \mid w_1^{k-1}) \simeq$$

$$P(w_1 | start) \cdot P(w_2 | w_1) \cdot P(w_3 | w_2) \cdots P(w_k | w_{k-1})$$

- We can think of the LM as a system that provides the probabilities $P(w_i|w_{i-1})$, which we then multiply.
 - o Or the probabilities $P(w_i|w_{i-2},w_{i-1})$ for a trigram LM.
 - o Or the probabilities $P(w_i|h_{i-1})$ for an LM that considers all the "history" (previous words) h_{i-1} , e.g., in an RNN LM.
- An LM typically provides a distribution P(w|h) showing how probable it is for every word $w \in V$ to be the next one.

RNN LM with GRU cells



 $\vec{o}_{i+1} = \operatorname{softmax}(W^{(o)}\vec{h}_i)$

New history. For $\vec{z}_i \approx \vec{0}$, same as \tilde{h}_i . For $\vec{z}_i \approx \vec{1}$, ignores \tilde{h}_i and maintains \vec{h}_{i-1} as \vec{h}_i .

GRU cell:

$$\tilde{h}_{i} = \tanh(\vec{r}_{i} \circ W^{(h)} \vec{h}_{i-1} + W^{(e)} \vec{e}_{i})$$

$$\vec{h}_{i} = \vec{z}_{i} \circ \vec{h}_{i-1} + (\vec{1} - \vec{z}_{i}) \circ \tilde{h}_{i}$$

$$\vec{r}_{i} = \sigma(W^{(r)} \vec{h}_{i-1} + U^{(r)} \vec{e}_{i})$$

$$\vec{z}_{i} = \sigma(W^{(z)} \vec{h}_{i-1} + U^{(z)} \vec{e}_{i})$$

 $E \int_{\vec{x}}$

 $\vec{e}_i = E\vec{x}_i$

 $W^{(h)}$

 $W^{(e)}$

Reset gate (σ is the sigmoid function).

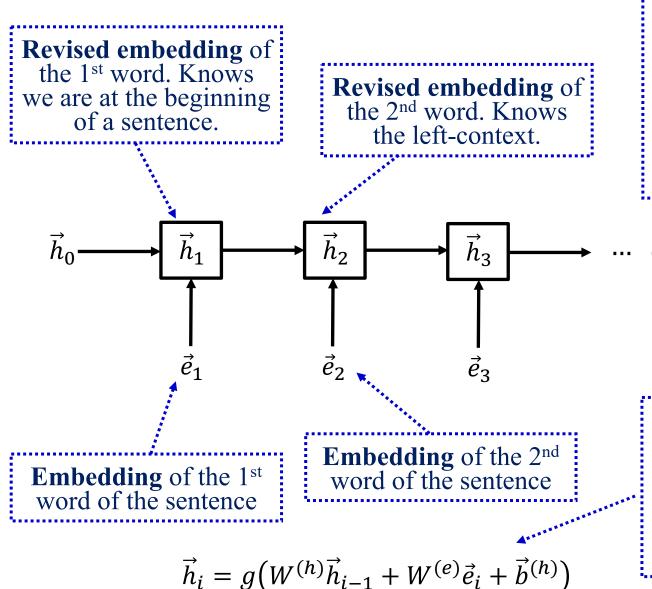
Update gate.

LSTM cells are similar, but with **more gates**. See http://colah.github.io/posts/201
5-08-Understanding-LSTMs/

More about RNNs

- Trained by backpropagation (with unrolled view).
 - o For each sentence (or window), feed it to the unrolled RNN, compute the loss and backpropagate, adding gradients obtained for the same matrix (e.g., same $W^{(h)}$ at each cell).
 - o GRU or LSTM cells help avoid vanishing gradients.
 - o The norms of the **gradients** can be **clipped** (when larger than a max value) to avoid **exploding gradients**.
 - Use layer normalization, not batch normalization in RNNs.
- We can also **learn** the **word embeddings** (*E*) with an RNN LM. Billions of **free training examples**!
 - We can re-use the word embeddings in other NLP tasks.
 - O With a large vocabulary, softmax is too slow (alternatives: small vocabulary, hierarchical softmax, negative sampling).

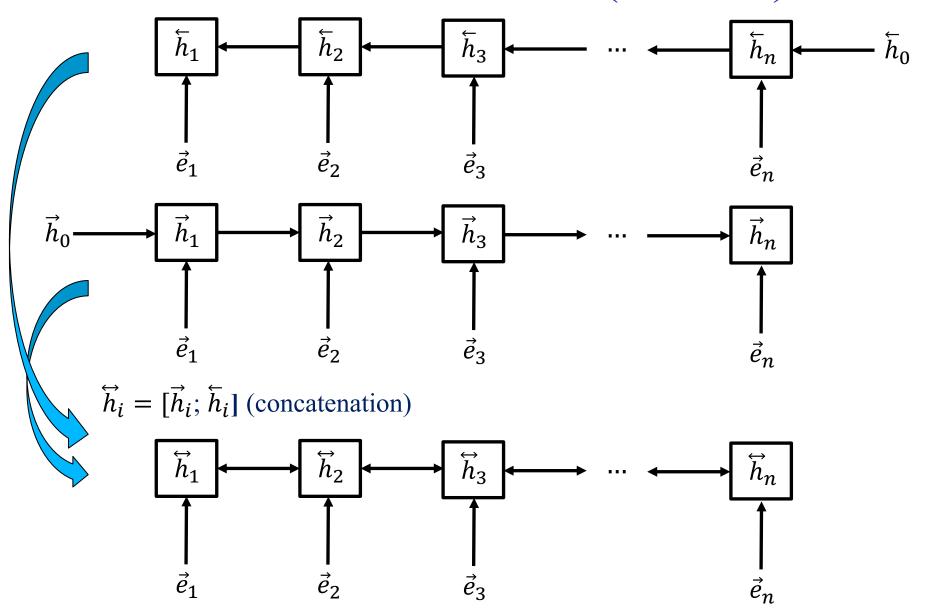
What about the right-context of each token?



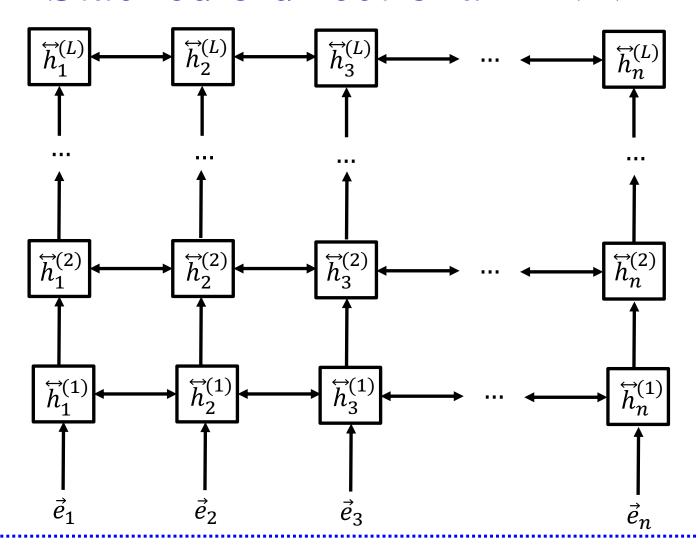
We can also treat the h_i vectors as the **memory** of the RNN, but in recent NLP work, it's easier to think of them as **revised** word embeddings.

g is an activation function (e.g., sigmoid). More complex update mechanisms in practice: LSTM or GRU cells.

Bidirectional RNN (biRNN)



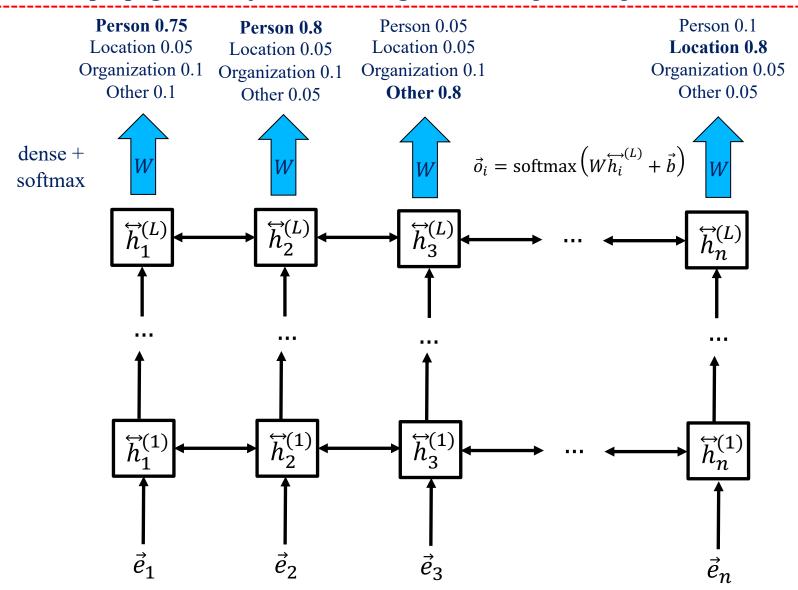
Stacked bidirectional RNN



Each layer revises the word embeddings of the previous (lower) layer. The embeddings become increasingly more context-aware and also increasingly more appropriate for the particular task we address...

Token classification with a stacked biRNN

Compare to the correct predictions (sum the cross-entropy loss for all token positions) and backpropagate to adjust all the weights, including the weights of the stacked biRNN.



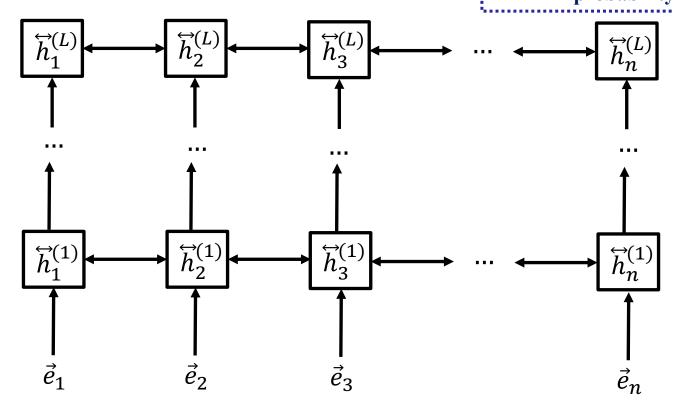
Text classification with stacked biRNN

Compare (via categorical cross entropy) the predicted \vec{o} to the correct 1-hot distribution and backpropagate to adjust all the weights, including the weights of the stacked biRNN.

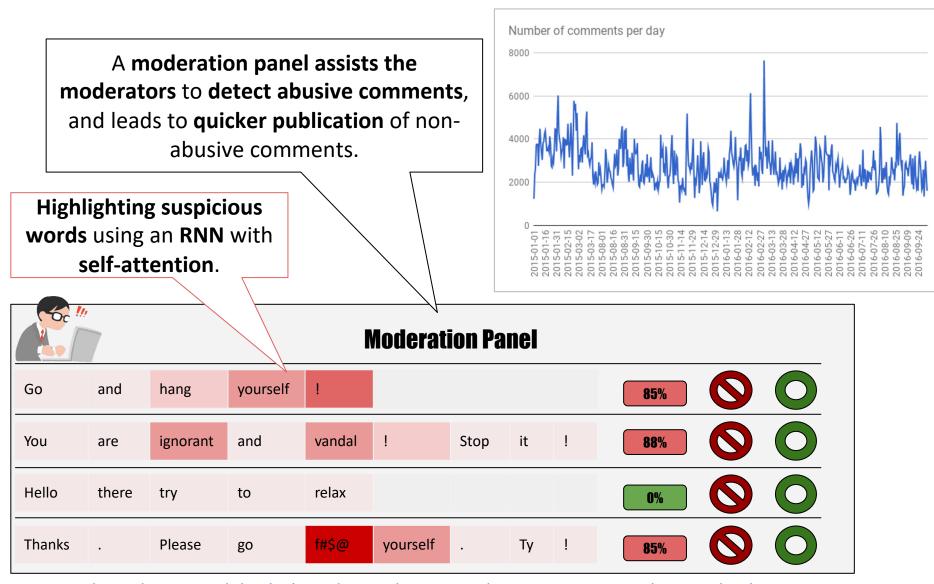
$$\vec{o} = \operatorname{softmax}(W\vec{h}_{max} + \vec{b})$$

$$\vec{h}_{max} = \left\langle \max\left(\overrightarrow{h}_{*,1}^{(L)}\right), \max\left(\overrightarrow{h}_{*,2}^{(L)}\right), \dots, \max\left(\overrightarrow{h}_{*,n}^{(L)}\right) \right\rangle^{T}$$

Global max-pooling creates a single vector containing the max per dimension of all the $\overleftrightarrow{h}_i^{(L)}$. We pass it through a dense layer and softmax (or MLP) to obtain a probability per class.

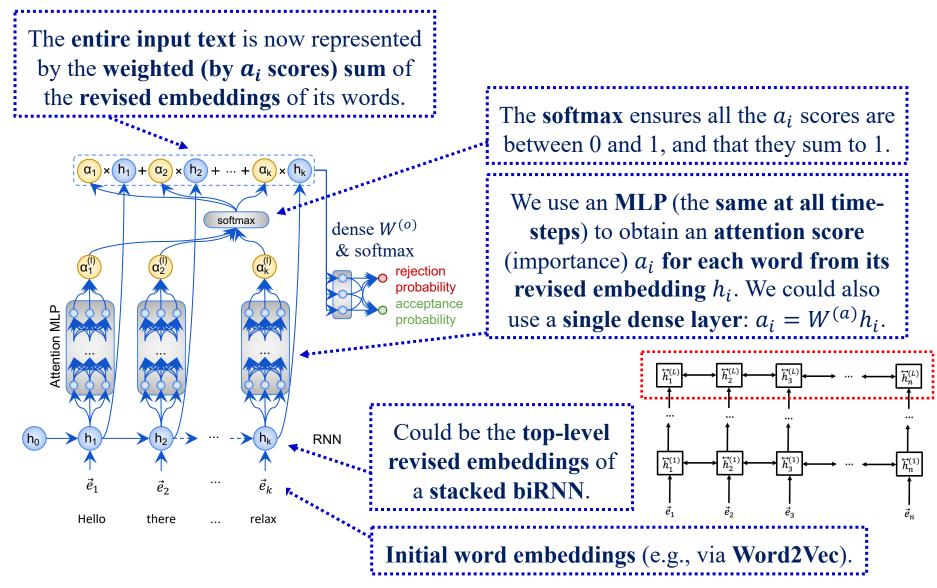


User comment moderation



J. Pavlopoulos, P. Malakasiotis and I. Androutsopoulos, "Deeper Attention to Abusive User Content Moderation", EMNLP 2017, http://nlp.cs.aueb.gr/pubs/emnlp2017.pdf.

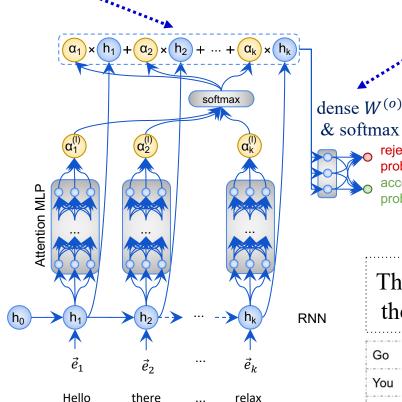
RNN with deep self-attention



J. Pavlopoulos, P. Malakasiotis and I. Androutsopoulos, "Deeper Attention to Abusive User Content Moderation", EMNLP 2017, http://nlp.cs.aueb.gr/pubs/emnlp2017.pdf.

RNN with deep self-attention

The entire input text is now represented by the weighted (by a_i scores) sum of the revised embeddings of its words.



We pass the weighted sum vector (point) through another dense layer and softmax to obtain a probability score for each class (here accept, reject).

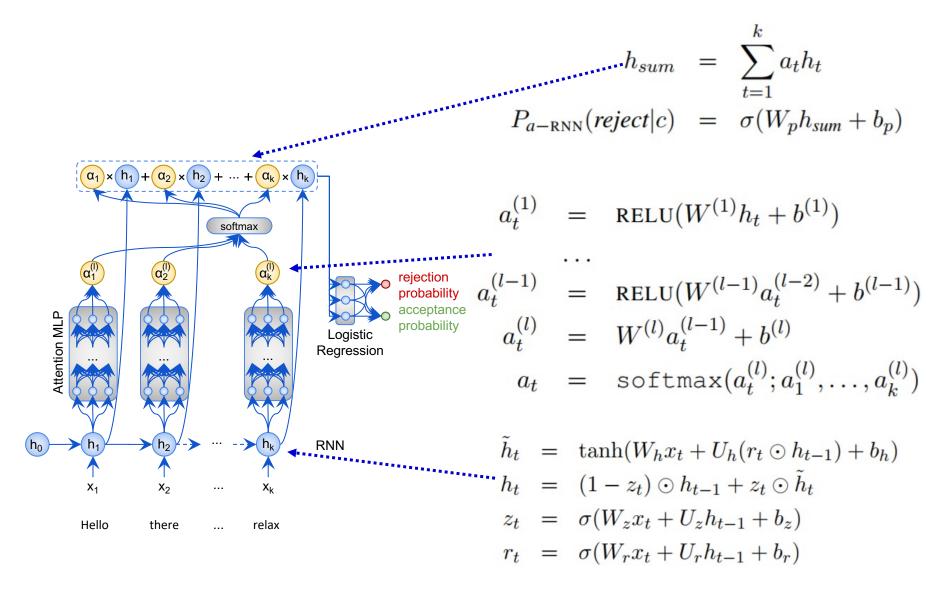
Compare to the correct predictions with a cross-entropy loss and backpropagate to adjust the weights of the entire neural net, including the MLP and RNN(s).

The attention scores a_i can also be used to highlight the words that influence the system's decision most.

Go	and	hang	yourself	!				
You	are	ignorant	and	vandal	!	Stop	it	!
Thanks		Please	go	W	yourself		ty	!

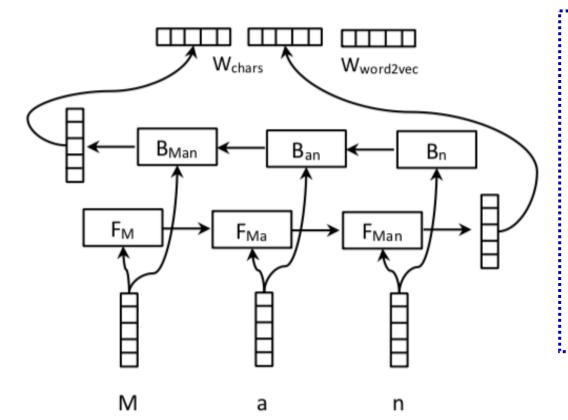
J. Pavlopoulos, P. Malakasiotis and I. Androutsopoulos, "Deeper Attention to Abusive User Content Moderation", EMNLP 2017, http://nlp.cs.aueb.gr/pubs/emnlp2017.pdf.

RNN with deep self-attention



J. Pavlopoulos, P. Malakasiotis and I. Androutsopoulos, "Deeper Attention to Abusive User Content Moderation", EMNLP 2017, http://nlp.cs.aueb.gr/pubs/emnlp2017.pdf.

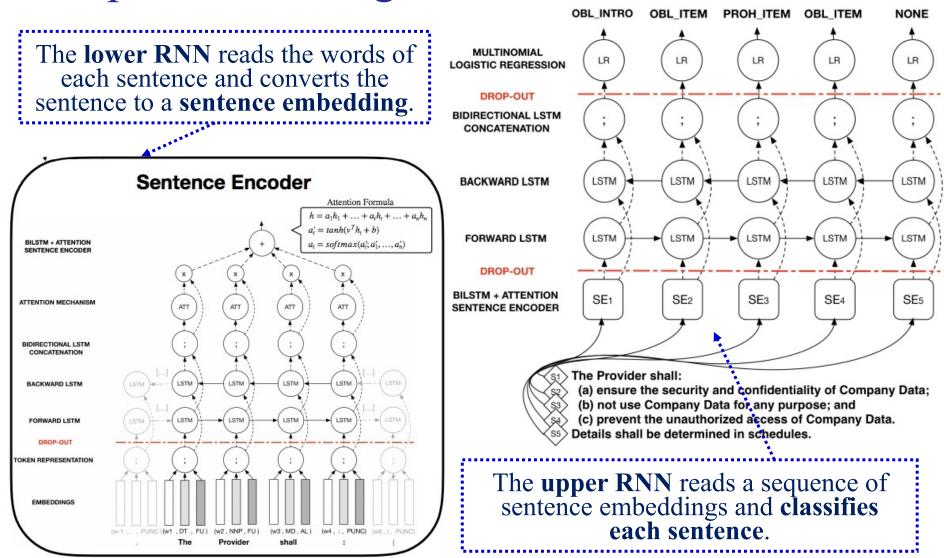
RNNs that produce word embeddings from character embeddings



Word embedding layer,
part of a larger network.
We concatenate the
word embedding we get
from the character-level
biLSTM with the
Word2Vec embedding.
The character
embeddings are learned
during back-propagation.

G, Bekoulis, J, Deleu, T, Demeester, C. Develder, "Joint entity recognition and relation extraction as a multi-head selection problem", Expert Systems with Applications, Vol, 114, pp. 34-45, 2018. Figure from the pre-print https://arxiv.org/abs/1804.07847.

Sequence labeling with a Hierarchical RNN



I. Chalkidis, I. Androutsopoulos, A. Michos, "Obligation and Prohibition Extraction Using Hierarchical RNNs", ACL 2018. http://www.aclweb.org/anthology/P18-2041

Legal judgment prediction for ECHR cases

Case ID: 001-148227 Violated Articles: Article 3 Predicted Violation: YES (0.97%)

- 1. The applicant was born in 1955 and lives in Kharkiv.
- 2. On 5 May 2004 the applicant was arrested by four police officers on suspicion of bribe taking.
 The police officers took him to the Kharkiv Dzerzhynskyy District Police Station, where he was held overnight.
 According to the applicant, the police officers beat him for several hours, forcing him to confess.
- 3. On 6 May 2004 the applicant was taken to the Kharkiv City Prosecutor's Office. He complained of ill-treatment to a senior prosecutor from the above office. The prosecutor referred the applicant for a forensic medical examination.
- 4. On 7 May 2004 the applicant was diagnosed with concussion and admitted to hospital.
- 5. On 8 May 2004 the applicant underwent a forensic medical examination, which established that he had numerous pruises on his face, chest, legs and arms, as well as a damaged tooth.
- 6. On 11 May 2004 criminal proceedings were instituted against the applicant on charges of bribe-taking. They were eventually terminated on 27 April 2007 for lack of corpus delicti.
- 7. On 2 June 2004 the applicant lodged another complaint of ill treatment with the Kharkiv City Prosecutor 's Office

Figure 1: Attention over words (colored words) and facts (vertical heat bars) as produced by HAN.

Words with high attention scores.

Sentences with **high attention** scores.

Biased against particular locations?

I. Chalkidis, I. Androutsopoulos and N. Aletras, "Neural Legal Judgment Prediction in English", ACL 2019. https://www.aclweb.org/anthology/P19-1424/ 21

RNNs for Machine Translation

From the slides of R. Socher's course "Deep Learning for NLP", 2015. http://cs224d.stanford.edu/

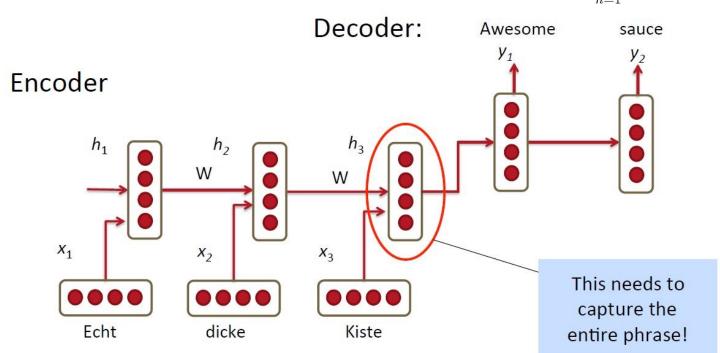
Encoder:
$$h_t = \phi(h_{t-1}, x_t) = f\left(W^{(hh)}h_{t-1} + W^{(hx)}x_t\right)$$

Decoder:
$$h_t = \phi(h_{t-1}) = f\left(W^{(hh)}h_{t-1}\right)$$

$$y_t = softmax\left(W^{(S)}h_t\right)$$

Minimize cross entropy error for all target words conditioned on source words

$$\max_{\theta} \frac{1}{N} \sum_{n=1}^{N} \log p_{\theta}(y^{(n)}|x^{(n)})$$



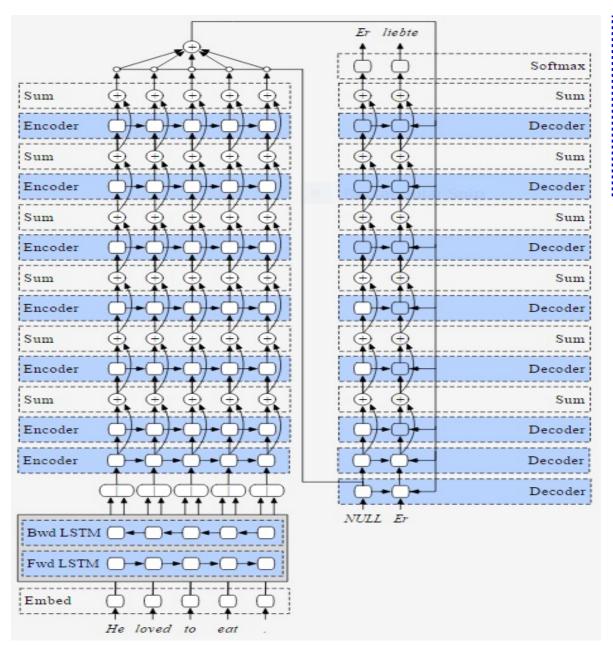
Different picture, same idea Embedding of the **previously** f = (La, croissance, économique, s'est, ralentie, ces, dernières, années, .)generated word. \mathbf{u}_i ecode Last \vec{h}_i of the \mathbf{p}_i encoder RNN. Treated as embedding of Recurrent the entire input sentence. Recurrent State From the slides of R. Word Representation Socher's course "Deep Continuous-space Learning for NLP", 2015. http://cs224d.stanford.edu/ 1-of-K coding

e = (Economic, growth, has, slowed, down, in, recent, years, .)

Kyunghyun Cho et al. 2014

4/22/15

RNN-based Machine Translation

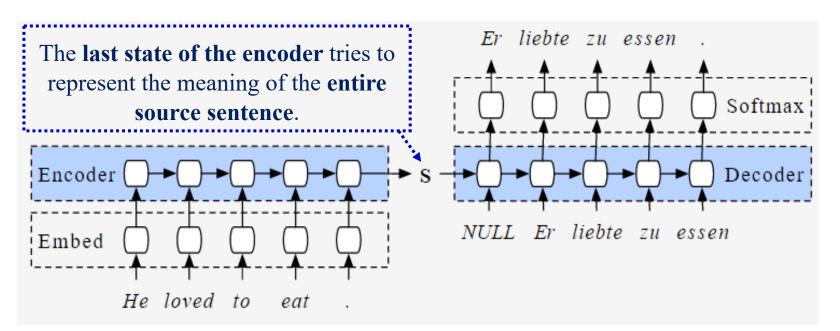


Google's paper: https://arxiv.org/abs/1609.08144

Images from Stephen Merity's http://smerity.com/articles/2016/google-nmt_arch.html

Easier to explain step by step...

Basic Encoder-Decoder NMT

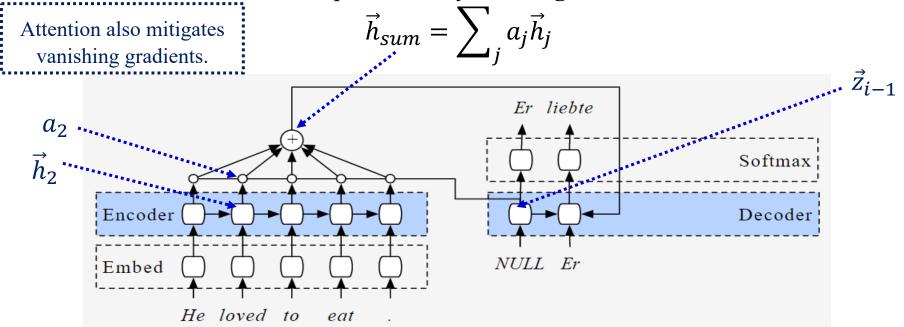


During training, at each decoding time-step, we can always use the correct previous word ("teacher forcing"); or we can randomly use the correct or (increasingly more often) the predicted previous word (scheduled sampling).

During testing (inference), we always use the **predicted previous word**; and we **greedily select the most probable next word**, or we **sample from the distribution of the next word**, or we use **beam search** to find the translation $y_1, ..., y_m$ with the highest probability: $p(y_1|z_1) p(y_2|y_1, z_2) p(y_3|y_2, z_3) ... p(y_m|y_{m-1}, z_m)$ where z_i are the states of the decoder.

Encoder-Decoder with attention

The source sentence is now represented by the weighted sum of the encoder states:



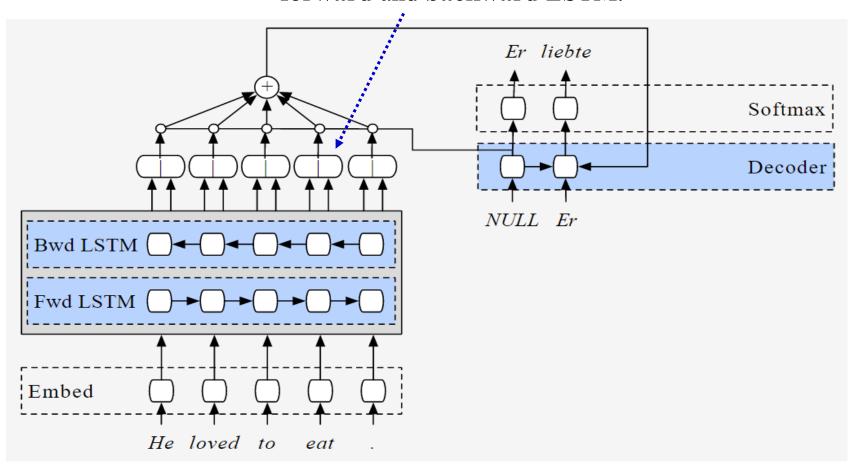
For each German word, the **attention scores** over the English words **change!** Each "**attention**" **weight** a_j is a **function** of the **corresponding encoder state** \vec{h}_j and the **previous state** \vec{z}_{i-1} **of the decoder** (memory of translation so far), e.g.: $\tilde{a}_j = v^T \cdot f(W^{(h)}\vec{h}_j + W^{(z)}\vec{z}_{i-1}) = v^T \cdot f(W[\vec{h}_j; \vec{z}_{i-1}]), \qquad a_j = softmax(\tilde{a}_j)$

with a **softmax** to make the a_i weights sum to 1.

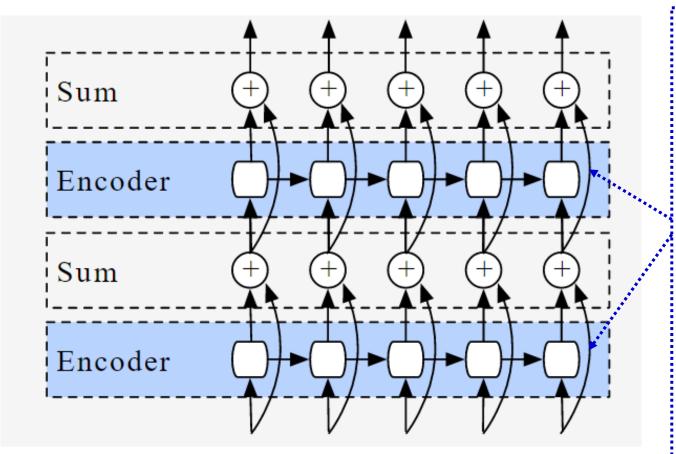
Google's paper: https://arxiv.org/abs/1609.08144
Images from Stephen Merity's https://smerity.com/articles/2016/google nmt arch.html

Bidirectional LSTM encoder

The encoder is now a **bidirectional LSTM**. The **encoder state** for the *j*-th word of the source sentence is the **concatenation** of the **corresponding states** of the **forward** and **backward** LSTM.

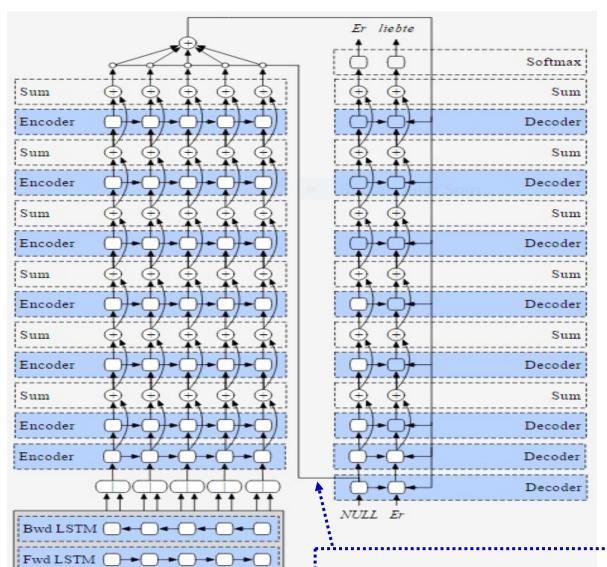


Stacking RNNs and residuals



"Residual" connections (a kind of skip-connections) helps fight vanishing gradients in backpropagation (sum-nodes copy the gradients to their inputs). Also allows upper layers to learn only modifications (differences) from representations of lower layers.

RNN-based Machine Translation



Embed

He loved to

Google's paper: https://arxiv.org/abs/1609.08144

Images from Stephen Merity's http://smerity.com/articles/2016/google nmt arch.html

Attention based on the previous state of the bottom decoder only, to speed up computations.

Additional optional reading slides.

Dropout vs. Variational Dropout

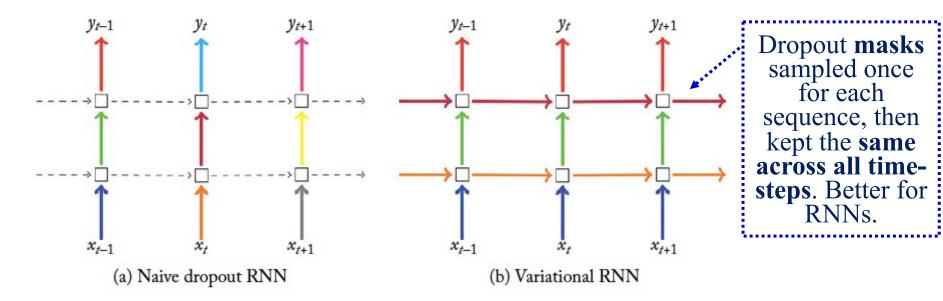
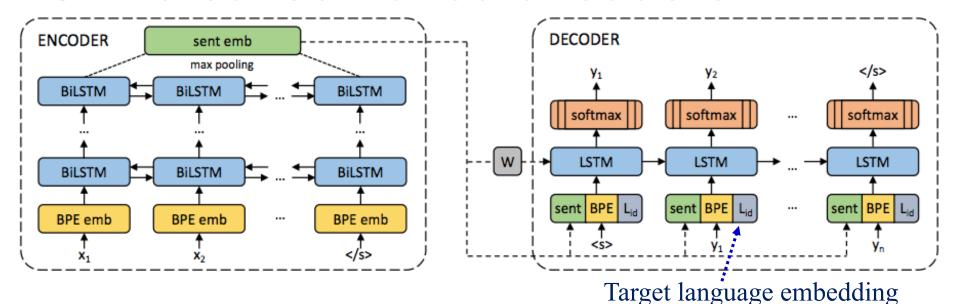


Figure 15.2: Gal's proposal for RNN dropout (b), vs. the previous suggestion by Pham et al. [2013], Zaremba et al. [2014] (a). Figure from Gal [2015], used with permission. Each square represents an RNN unit, with horizontal arrows representing time dependence (recurrent connections). Vertical arrows represent the input and output to each RNN unit. Colored connections represent dropped-out inputs, with different colors corresponding to different dropout masks. Dashed lines correspond to standard connections with no dropout. Previous techniques (naive dropout, left) use different masks at different time steps, with no dropout on the recurrent layers. Gal's proposed technique (Variational RNN, right) uses the same dropout mask at each time step, including the recurrent layers.

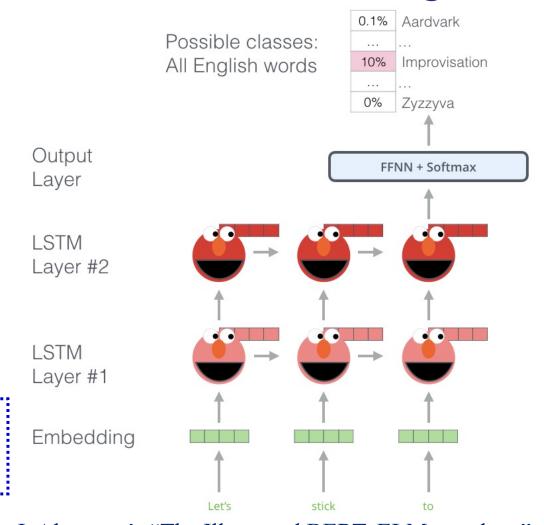
Figure from: Y. Goldberg, *Neural Network Models for Natural Language Processing*, Morgan & Claypool Publishers, 2017. See also https://adriangcoder.medium.com/a-review-of-dropout-as-applied-to-rnns-72e79ecd5b7b

Universal sentence encoders



- Laser: Trained on parallel corpora of 93 languages.
 - o Using the same encoder and decoder for all languages.
 - o Shared vocabulary of sub-word units (BPEs).
- E.g., we can **train a classifier** on **English tweets**, and use the **same trained classifier** on **Greek tweets**.
 - o In **both languages**, we use the **same encoder** to convert each **tweet** to a **feature vector**.
- M. Artetxe and H. Schwenk, "Massively Multilingual Sentence Embeddings for Zero-Shot Cross-Lingual Transfer and Beyond". https://code.fb.com/ai-research/laser-multilingual-sentence-embeddings/

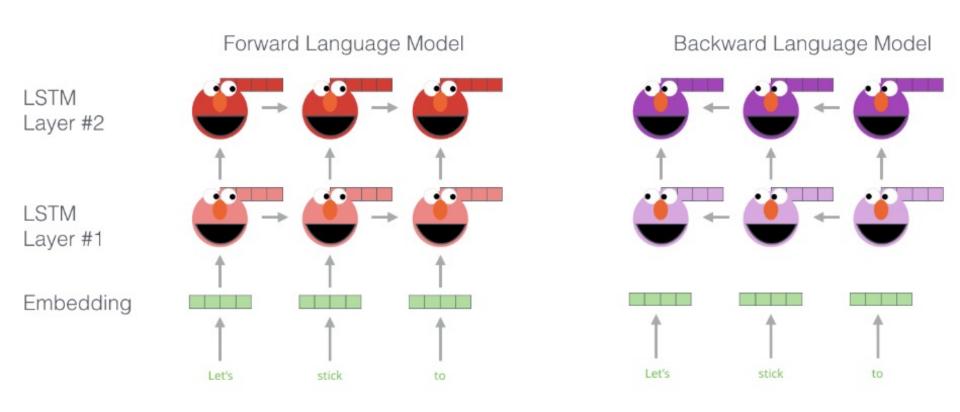
ELMo – Pretraining LMs to obtain context aware word embeddings



Bottom word embeddings generated by CNNs operating on characters.

Figures from J. Alammar's "The Illustrated BERT, ELMo, and co." http://jalammar.github.io/illustrated-bert/. ELMo paper: Peters et al. "Deep Contextualized Word Representations", NAACL-HLT 2018. http://aclweb.org/anthology/N18-1202

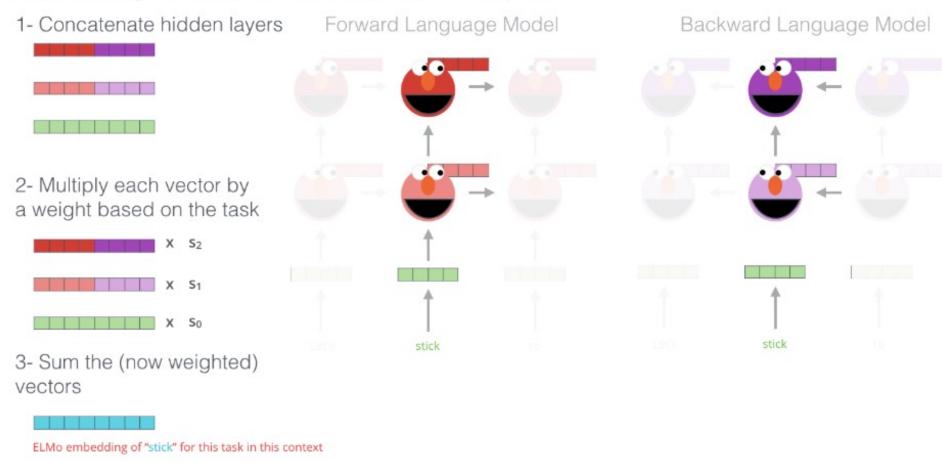
ELMo – Pretraining LMs to obtain context aware word embeddings



Figures from J. Alammar's "The Illustrated BERT, ELMo, and co." http://jalammar.github.io/illustrated-bert/. ELMo paper: Peters et al. "Deep Contextualized Word Representations", NAACL-HLT 2018. http://aclweb.org/anthology/N18-1202

ELMo – Pretraining LMs to obtain context aware word embeddings

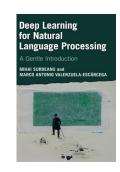
Embedding of "stick" in "Let's stick to" - Step #2

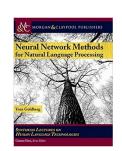


Figures from J. Alammar's "The Illustrated BERT, ELMo, and co." http://jalammar.github.io/illustrated-bert/. ELMo paper: Peters et al. "Deep Contextualized Word Representations", NAACL-HLT 2018. http://aclweb.org/anthology/N18-1202

Recommended reading

- M. Surdeanu and M.A. Valenzuela-Escarcega, *Deep Learning for Natural Language Processing: A Gentle Introduction*, Cambridge Univ. Press, 2024.
 - Chapters 11, 12, 14. See https://clulab.org/gentlenlp/text.html
 - Also available at AUEB's library.
- Y. Goldberg, Neural Network Models for Natural Language Processing, Morgan & Claypool Publishers, 2017.
 - o Mostly chapters 14–17.
- Jurafsky and Martin's, *Speech and Language Processing* is being revised (3rd edition) to include DL methods.
 - http://web.stanford.edu/~jurafsky/slp3/







Recommended reading

- A. Zhang et al., *Dive into Deep Learning*.
 - Freely available at: https://d21.ai/
 - See Chapters 9 and 10 for RNNs.

