# Sequence 2 Sequence

L101: Machine Learning for Language Processing Andreas Vlachos



## Structured prediction reminder

Given an input x (e.g. a sentence) predict y (e.g. a PoS tag sequence, cf lecture 5):

$$\hat{y} = rg \max_{y \in \mathcal{Y}} score(x,y)$$

Where Y is rather large and often depends on the input (e.g.  $L^{|x|}$  in PoS tagging)

#### Various approaches:

- Linear models (structured perceptron)
- Probabilistic linear models (conditional random fields)
- Generative models (hidden Markov models)

#### Most common structures

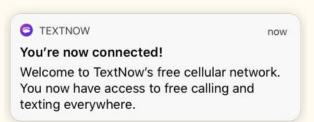
As input?

#### SEVINDIRICI HABER YUHANNA

1 Her şeyin başlangıcından önce Tannsal Söz vardı. Tannsal Söz Tann'yla birlikteydi ve Tann neyse Tannsal Söz O'ydu. <sup>2</sup>Başlangıçta Tann'yla birlikteydi. <sup>3</sup>Tann her şeyi O'nun aracılığıyla oluşturdu ve olanlardan hiçbiri O'nsuz olmadı. <sup>4</sup>Her varlığa yaşam veren O'ydu ve yaşam insanların Işığı'ydı\*. <sup>5</sup>Işık karanlığı aydınlatır, karanlık onu alt edemez.

<sup>6</sup>Tanın tarafından gönderilmiş bir adam ortaya çıktı, adı Yahya idi. <sup>7</sup>Tanıklık etmeye geldi. Işık için tanıklık etsin ve herkes onun aracılığıyla imana gelsin diye geldi. <sup>8</sup>Kendisi o Işık değildi, sadece Işık için tanıklık etmeye geldi. Dies ist ein Blindtext. An ihm lässt sich vieles über die Schrift ablesen, in der er gesetzt ist. Auf den ersten Blick wird der Grauwert der Schriftfläche sichtbar. Dann kann man prüfen, wie gut die Schrift zu lesen ist und wie sie auf den Leser wirkt. Dies ist ein Blindtext. An ihm lässt sich vieles über die Schrift ablesen, in der er gesetzt ist. Auf den ersten Blick wird der Grauwert der Schriftfläche sichtbar. Dann kann man prüfen, wie gut die Schrift zu lesen ist und wie sie auf den Leser wirkt.

As output?



Someone who came in contact with you tested positive or has shown symptoms for COVID-19 & recommends you self-isolate/get tested. More at

Natural language, i.e. sequences of character, words, sentences!

Today: focus on language-to-language methods, a.k.a. seq2seq, encoder-decoder

## Language modelling

How likely is that a sequence of words comes from a particular language (e.g. English)?

Odd sounding problem. Applications:

- speech recognition
- machine translation
- grammatical error detection
- etc.

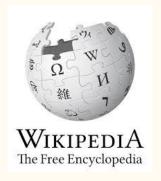
### Problem setup

We want to learn a model that gives us:

$$P(\mathbf{x}), \mathbf{for} \ orall \mathbf{x} \in V^{maxN}$$

Training?

As much text as we can get!







## Language modelling

Decompose the probability of the sentence x into conditional probabilities of each word given the previous ones:

$$P(\mathbf{x}) = P(x_1)P(x_2|x_1)P(x_3|x_2,x_1)\dots P(x_N|x_1,\dots,x_{N-1})$$

These are typically (<u>until 2010</u>) estimated using word counts:

$$P(x_n|x_{n-1...x_1}) = rac{counts(x_1...x_{n-1},x_n)}{counts(x_1...x_{n-1})}$$

Any problems?

#### **Sparsity!** Solutions:

- Markov assumption, i.e. N-gram language models
- smoothing: interpolation between models, Kneser-Ney, stupid back-off, etc.

### What is a language model?

A giant logistic regression classifier over words:

$$p(x_n = k | x_{n-1} \dots x_1) = rac{\exp(\mathbf{w}_k \cdot \phi(x_{n-1} \dots x_1))}{\sum_{k'=1}^{|\mathcal{V}|} \exp(\mathbf{w}_{k'} \cdot \phi(x_{n-1} \dots x_1))} = softmax(\mathbf{W} \cdot \phi(x_{n-1} \dots x_1))$$

But terribly inefficient using counts/one-hot encoding as features!

#### A language generator:

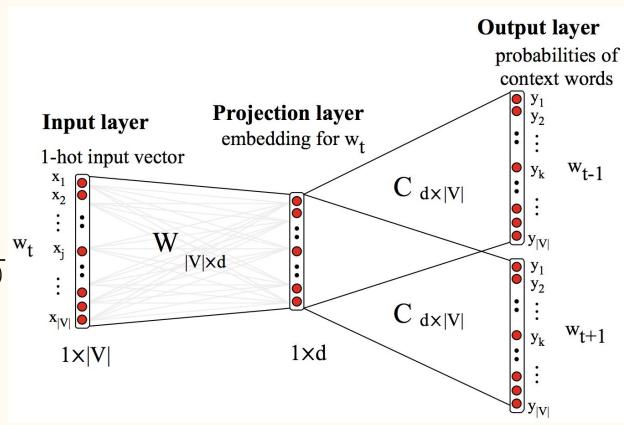
- Sample repeatedly from p(x), each time adding the words in the context
- Stop when the <END> of the sentence token is sampled

## Skipgram word embeddings

Skipgram (<u>Mikolov et al.</u> 2013) is a giant word-given-word classifier with learned features:

$$P(w_{t-1}|w_t) = rac{\exp(\mathbf{c_{t-1}}\cdot\mathbf{w_t})}{\sum_{c'\in V}\exp(\mathbf{c'}\cdot\mathbf{w_t})}^{\mathbf{w_t}}$$

(each word has two embeddings)

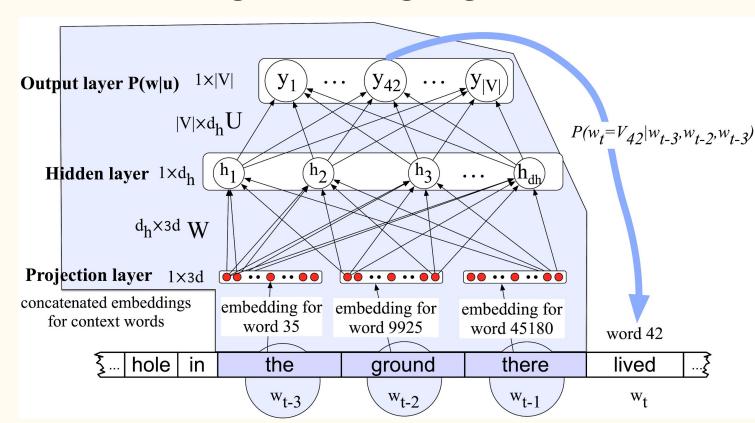


### A Feedforward NN N-gram Language model

Using 3 words as context instead of 1. Limitations?

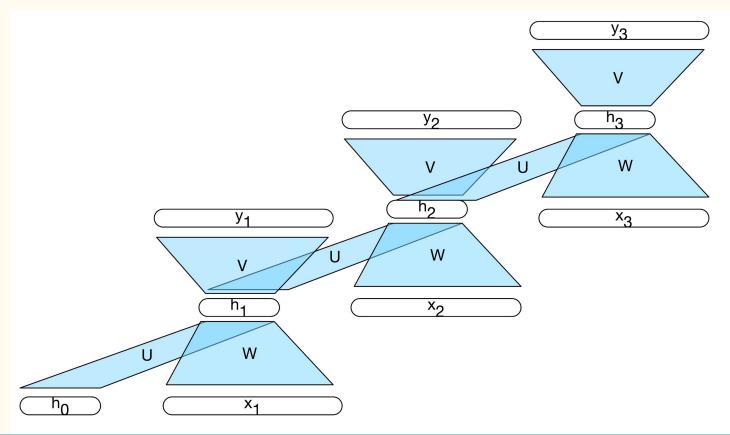
Context used still limited

NN has to learn how to use the same word in each context place separately



#### Recurrence

Each word is processed using the same weight matrices.



#### Recurrent Neural Network

$$egin{aligned} p(x_n|x_{n-1}\dots x_1) &= softmax(\mathbf{V}\cdot h_n) \ &= softmax(\mathbf{V}\cdot \phi(x_{n-1}\dots x_1)) \ h_n &= g(\mathbf{U}\cdot h_{n-1} + \mathbf{W}\cdot x_n) \end{aligned}$$

- V is the output layer, like the weights of logistic regression
- W is the word embeddings dictionary (can be pre-trained/fine-tuned)
- g is a nonlinear function, e.g. tanh
- U determines how to use the representation of the context  $h_{t-1}$

## Backpropagation through time

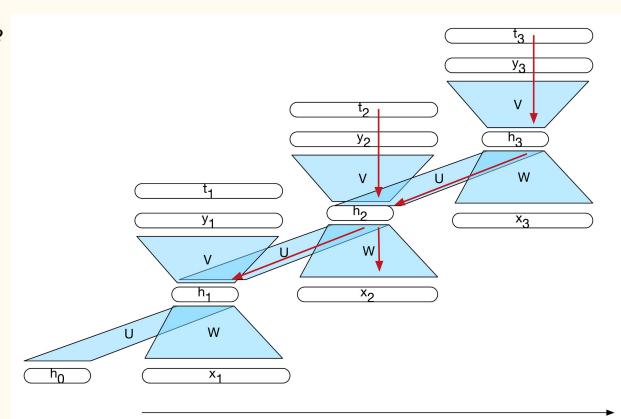
Why not simple backpropop?

Each  $\mathbf{h}_{t}$  affects  $\mathbf{y}_{t}$  and every  $\mathbf{y}_{t'>t}$  afterwards

The loss calculation needs to take all into account

Unroll the network for a fixed number of steps

Still uses unlimited context during testing

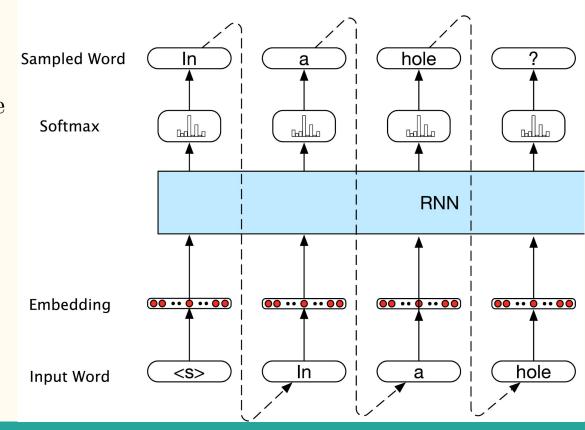


### Decoding, i.e. generation

- 1. Sample a word
- 2. Feed its embedding
- 3. Repeat until end of sentence

Many options for decoding

- max
- random
- top-k
- <u>nucleus</u> (Holtzman et al.)
- <u>temperature scaling</u>
- beam search (next lecture)



## Typical combinations

Sequence encoder and classifier

Seq2seq, encoder-decoder, language generator conditioned on sequence

one to one one to many many to one many to many many to many

Which tasks fit which variant

- PoS tagging?
- Machine translation?
- Image captioning?

Conditional language generator

http://karpathy.github.io/2015/05/21/rnn-ctiveness/

Sequence encoder, token tagger

## A few more pointers

- Bidirectional RNN encoders are commonly used
  - Sentence representations avoiding being biased by the last tokens
  - Word representations taking into account context left and right
- Multiple hidden layers are commonly used (stacked RNNs)
  - More demanding computationally
  - Able to learn more high level features
- Long-Short Term Memory Networks (Hochreiter and Schmidhuber, 1997) were introduced to handle long-range dependencies
- <u>Convolutional Neural Networks</u> are also used relying on multiple layers to mitigate the effect of the fixed window
- While using words as the modelling unit, using characters and <u>subwords</u> helps deal with rare/unknown words

## Long-range dependencies

RNNs don't handle them well, but Long-Short Term Memory Networks (Hochreiter and Schmidhuber, 1997) do:

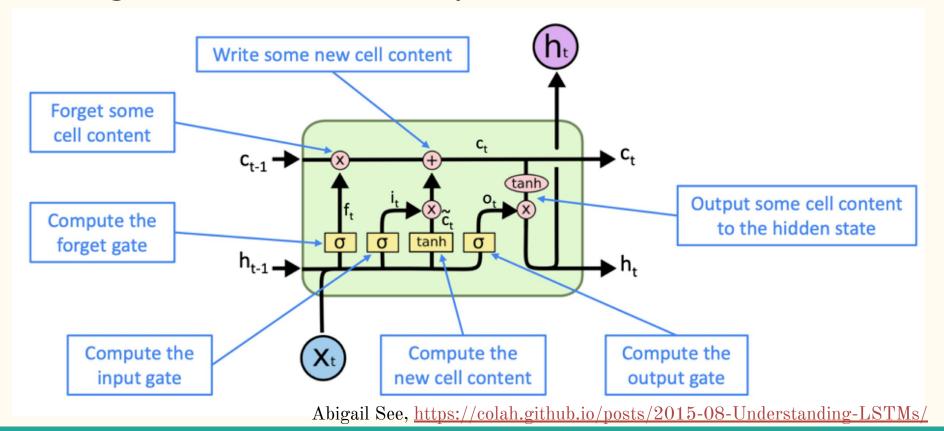
- introduce a memory cell, running parallel to the hidden state
- forget gate that decides which part of the memory to drop
- input gate that decides which part of the input to add to the memory
- output gate that decides which part of the memory to use in the hidden state

#### Advantages:

- Memory cell allows us to keep information not immediately needed
- Addresses the issue of vanishing/exploding gradients
  - see gradient clipping as an alternative

Main disadvantage: more parameters to learn, but usually worth it

### Long-short term memory networks

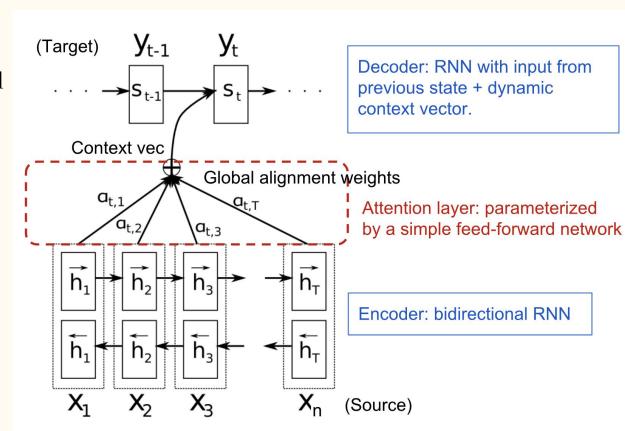


#### Attention

Generating a sentence based on one vector suboptimal: not all inputs relevant to every output

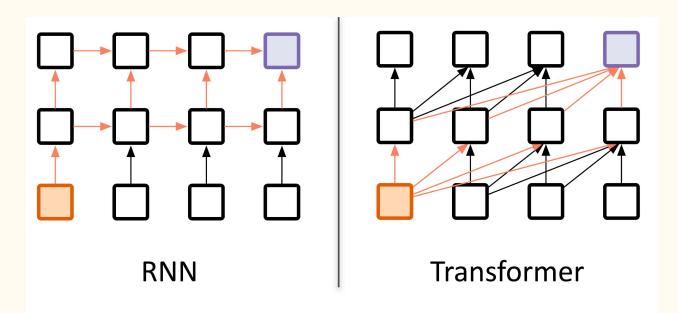
Allow the model to use hidden states of the input apart from the last one

Attention intuitively works as an alignment mechanism, but it is not clear <a href="how/why">how/why</a>



https://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html

#### Self-attention instead of recurrence

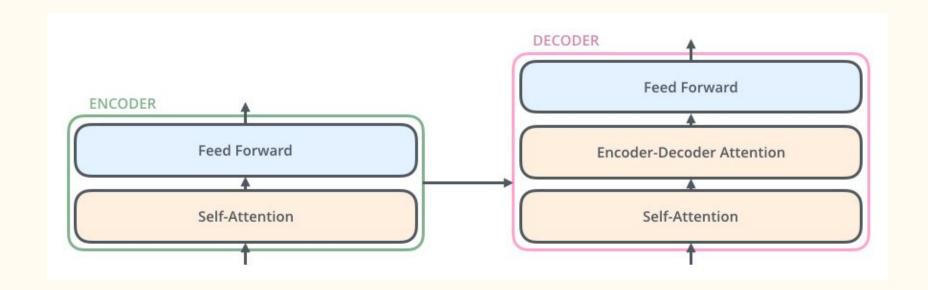


Arriana Bisazza (AthNLP 2019)

Key idea behind the Transformers (<u>Vaswani et al. 2017</u>)

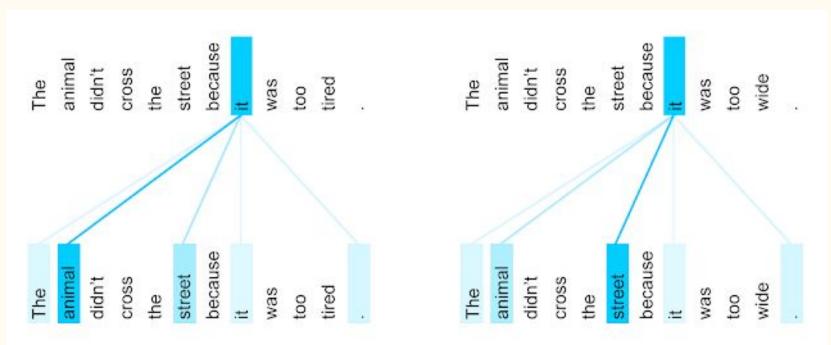
- Better parallelization, i.e. faster, more data, etc.
- Can be seen as a fully connected <u>graph neural network</u>

#### Transformers - overview



http://jalammar.github.io/illustrated-transformer/

### Transformers - multiple layers of self-attention



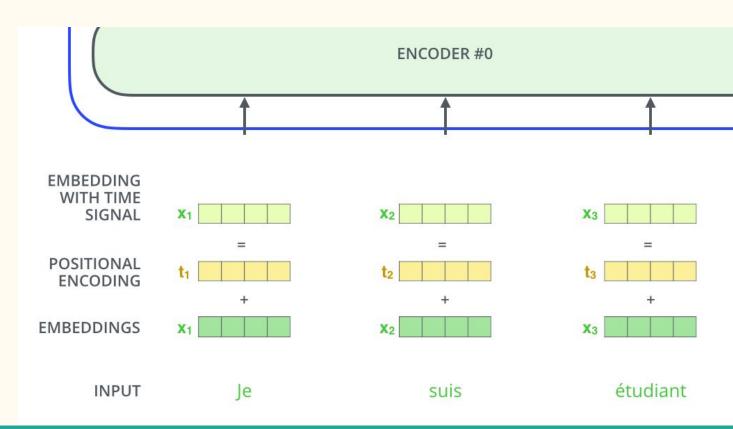
Each layer has multiple attention heads which <u>can be pruned after training</u>

 $\frac{https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html}{}$ 

## Transformers - positional encoding

No recurrence, but order is taken into account via the positional embeddings

http://jalammar.github.io/illustrated-transformer/



#### Transformers - decoder

http://jalammar.git

hub.io/illustrated-t

ransformer/

Linear + Softmax Kencdec Vencdec Same as encoding but can only **ENCODERS DECODERS** self-attend to tokens already generated and **EMBEDDING** WITH TIME attends to the SIGNAL input too **EMBEDDINGS** 

étudiant

suis

le

INPUT

**PREVIOUS** 

**OUTPUTS** 

## Encoder only: BERT (Devlin et al., 2018)

#### BERT: <u>Bidirectional Encoder Representations from Transformers</u>, i.e.:

- take the transformer encoder stack (self-attention, positional encodings, etc.)
- download a lot of text (sub-word tokenized)
- add special tokens for the sentence beginning and separators
- train two models: left-to-right and right-to-left using two objectives:
  - Masked language modelling: predict words missing at random from the text
  - Next sentence prediction: predict whether the next sentence was the one in the text or not

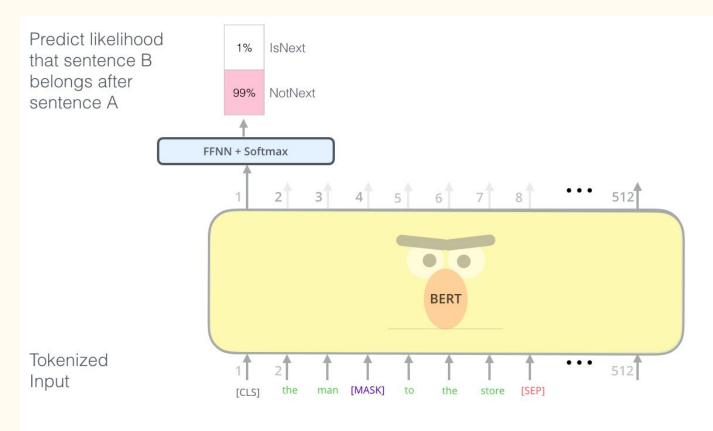
#### Typically considered the baseline method to beat:

- use it pre-trained as an input encoder/feature extractor
- fine-tune it to produce token/sentence embeddings for the task

#### BERT

Task 1: masked language modelling

Task 2: next sentence prediction



http://jalammar.github.
io/illustrated-bert/

Input

[CLS] the man [MASK] to the store [SEP] penguin [MASK] are flightless birds [SEP]

Sentence A Sentence B

#### Decoders as encoders?

Yes! Generalised pre-training (GPT)!

Just take the embedding produced by the decoder before predicting the word

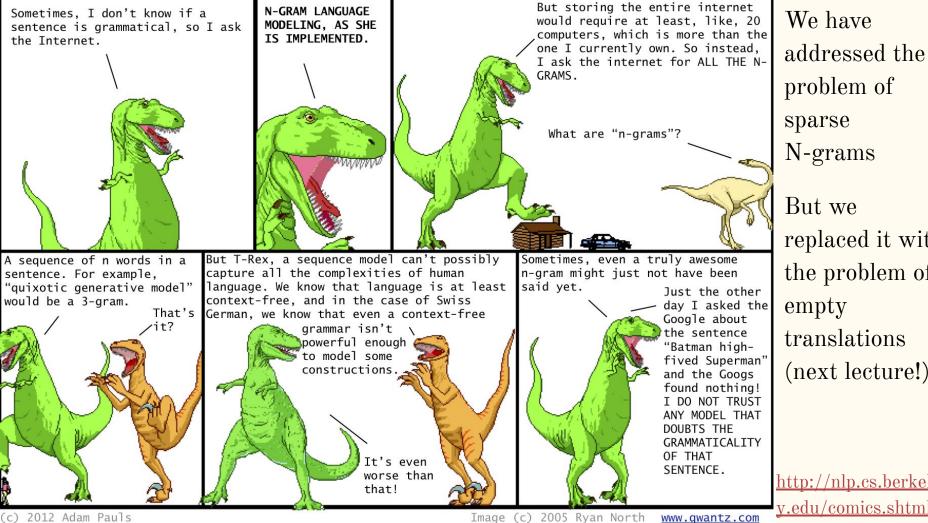
Simpler architecture, better scaling!

Decoders is all we need?

- An overkill when we just need an embedding (millions vs billions of params)
- Dedicated encoders can produce superior representations (Qorib et al. 2024)
- With enough labeled data encoder only has better accuracy (<u>Wang et al. 2024</u>)

## Bibliography

- Jurafsky and Martin chapters on <u>neural language models</u> and <u>RNNs</u> (RNN figures are from there unless otherwise stated)
- This <u>blog</u> explains RNNs and BPTT with code
- The <u>deep learning book</u>, chapter 10
- Two blogs about transformers and J&M chapter
- <u>BERT survey</u>
- Noah Smith's <u>introduction to contextual word embeddings</u>
- Seq2Seq can be used for generating sequences that are not words (only) such as semantic parsing



problem of sparse N-grams But we

replaced it with the problem of empty translations (next lecture!)

http://nlp.cs.berkele v.edu/comics.shtml