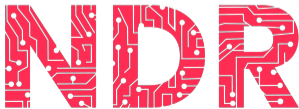


An Introduction to Neural Machine Translation

Prof. John D. Kelleher
@johndkelleher

ADAPT Centre for Digital Content Technology
Dublin Institute of Technology, Ireland

June 25, 2018



Outline

The Neural Machine Translation Revolution

Neural Networks 101

Word Embeddings

Language Models

Neural Language Models

Neural Machine Translation

Beyond NMT: Image Annotation

Microsoft Translator launching Neural Network based translations for all its speech languages

Rate this article



Microsoft Translator November 15, 2016

Share 461

287

in 0

0

Microsoft Translator is now powering all speech translation through state-of-the-art neural networks.

All speech translation apps that use this service, such as Skype Translator and the Microsoft Translator app for mobile devices, are now using neural network technology. Furthermore, the technology is available to all developers and end-users who want to use the Microsoft Translator speech API to integrate the technology into their favorite apps and services.

In addition to the nine languages supported by the Microsoft Translator speech API, namely Arabic, Chinese Mandarin, English, French, German, Italian, Brazilian Portuguese, Russian and Spanish, neural networks also power Japanese and Korean text translations. These eleven languages together represent more than 80% of the translations performed daily by Microsoft Translator.

Neural network technology has been used for the last few years in many artificial intelligence scenarios, such as speech and image processing. Many of these capabilities are available through [Microsoft Cognitive services](#). Neural networks are making in-roads into the machine translation industry, providing major advances in translation quality over the existing industry-standard Statistical Machine Translation (SMT) technology. Because of how the technology functions, neural networks better capture the context of full sentences before translating them, providing much higher quality and more human-sounding output.

Higher quality neural translations for a bunch more languages



Barak Turovsky
Product Lead (and proud
Russian speaker!), Google
Translate

Published Mar 6, 2017

Last November, people from Brazil to Turkey to Japan discovered that Google Translate for their language was suddenly more accurate and easier to understand. That's because we introduced [neural machine translation](#)—using deep neural networks to translate entire sentences, rather than just phrases—for eight languages overall. Over the next couple of weeks, these improvements are coming to Google Translate in many more languages, starting right now with Hindi, Russian and Vietnamese.

Neural translation is a lot better than our previous technology, because we translate whole sentences at a time, instead of pieces of a sentence. (Of course

Artificial intelligence

Facebook

applied mathematics

artificial intelligence

artificial neural networks

Facebook finishes its move to neural machine translation

Posted Aug 3, 2017 by [John Mannes \(@JohnMannes\)](#)



Next Story



Facebook announced [this morning](#) that it had completed its move to neural machine translation —

Image from

<https://techcrunch.com/2017/08/03/facebook-finishes-its-move-to-neural-machine-translation/>

Linguee's Founder Launches DeepL in Attempt to Challenge Google Translate

by [Florian Faes](#) on August 30, 2017



Barely two years after bursting into the translation tech scene, [neural machine translation](#) (NMT) is everything the MT community is talking about. Microsoft, Google, [Facebook](#), and other large technology companies have all transitioned to NMT, as did the [European Patent Office](#) and the [World Intellectual Property Organization](#). Even end-buyers are starting to build their own systems based on [open-source models](#).

Image from <https://slator.com/technology/linguees-founder-launches-deepl-attempt-challenge-google-translate/>

[//slator.com/technology/linguees-founder-launches-deepl-attempt-challenge-google-translate/](https://slator.com/technology/linguees-founder-launches-deepl-attempt-challenge-google-translate/)

Neural Networks 101

What is a function?

A **function** maps a set of inputs (numbers) to an output (number)¹

$$\textit{sum}(2, 5, 4) \rightarrow 11$$

¹This introduction to neural network and machine translation is based on: Kelleher (2016)

What is a WEIGHTEDSUM function?

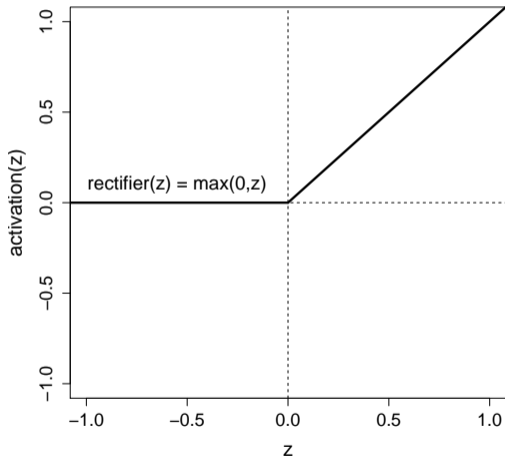
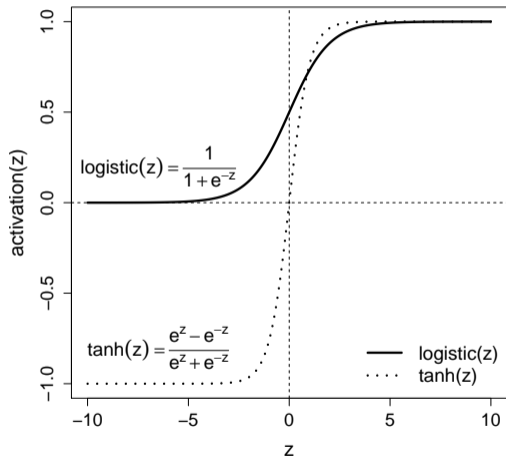
$$\begin{aligned} \text{WEIGHTEDSUM}(\underbrace{[x_1, x_2, \dots, x_m]}_{\text{Input Numbers}}, \underbrace{[w_1, w_2, \dots, w_m]}_{\text{Weights}}) \\ = (x_1 \times w_1) + (x_2 \times w_2) + \dots + (x_m \times w_m) \end{aligned}$$

$$\begin{aligned} \text{WEIGHTEDSUM}([3, 9], [-3, 1]) \\ = (3 \times -3) + (9 \times 1) \\ = -9 + 9 \\ = 0 \end{aligned}$$

What is an ACTIVATION function?

An ACTIVATION function takes the output of our WEIGHTEDSUM function and applies another mapping to it.

What is an ACTIVATION function?



What is an ACTIVATION function?

ACTIVATION =

$$\text{LOGISTIC}(\underbrace{\text{WEIGHTEDSUM}([x_1, x_2, \dots, x_m])}_{\text{Input Numbers}}, \underbrace{[w_1, w_2, \dots, w_m]}_{\text{Weights}})$$

$$\begin{aligned}\text{LOGISTIC}(\text{WEIGHTEDSUM}([3, 9], [-3, 1])) \\ &= \text{LOGISTIC}((3 \times -3) + (9 \times 1)) \\ &= \text{LOGISTIC}(-9 + 9) \\ &= \text{LOGISTIC}(0) \\ &= 0.5\end{aligned}$$

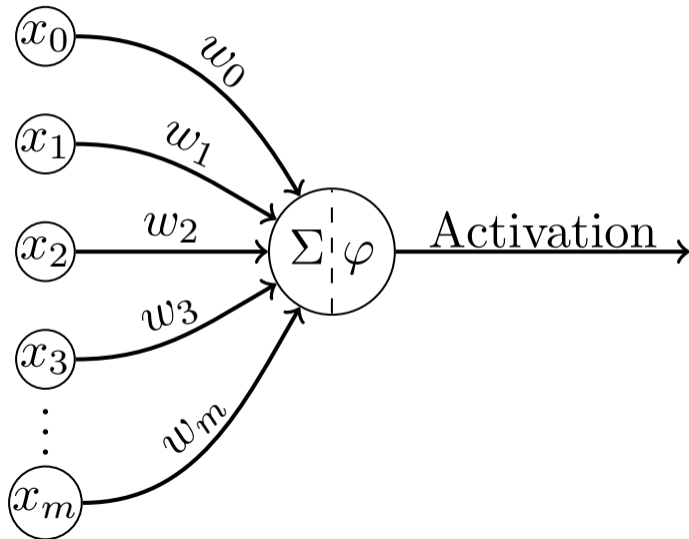
What is a NEURON?

The simple list of operations that we have just described defines the fundamental building block of a neural network: the NEURON.

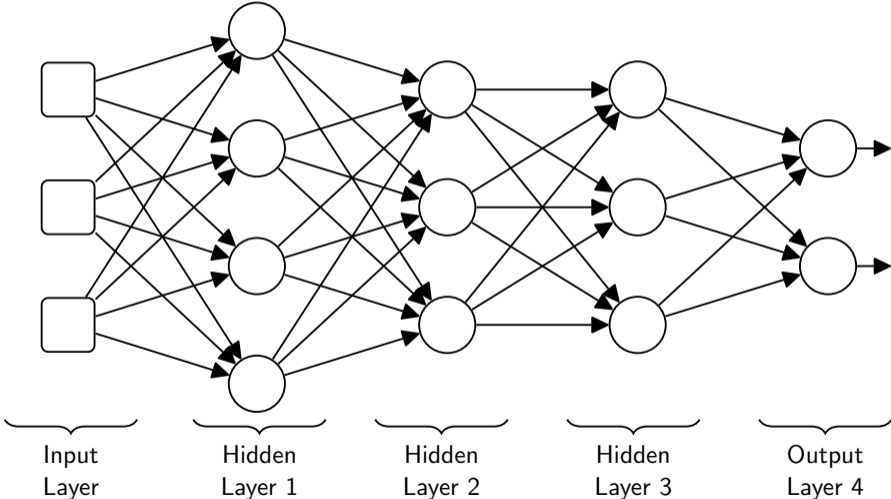
NEURON =

ACTIVATION(WEIGHTEDSUM($\underbrace{[x_1, x_2, \dots, x_m]}_{\text{Input Numbers}}, \underbrace{[w_1, w_2, \dots, w_m]}_{\text{Weights}}$))

What is a NEURON?



What is a NEURAL NETWORK?



Training a NEURAL NETWORK

- ▶ We train a neural network by iteratively updating the weights
- ▶ We start by randomly assigning weights to each edge
- ▶ We then show the network examples of inputs and expected outputs and update the weights using BACKPROPAGATION so that the network outputs match the expected outputs
- ▶ We keep updating the weights until the network is working the way we want

Word Embeddings

Word Embeddings

- ▶ Language is sequential and has **lots of words**.

“a word is characterized by the company it keeps”

— Firth, 1957

Word Embeddings

1. Train a network to predict the word that is missing from the middle of an n-gram (or predict the n-gram from the word)
2. Use the trained network weights to represent the word in vector space.

Word Embeddings

Each word is represented by a vector of numbers that positions the word in a multi-dimensional space, e.g.:

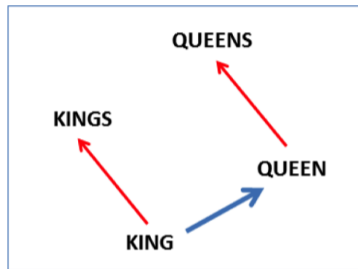
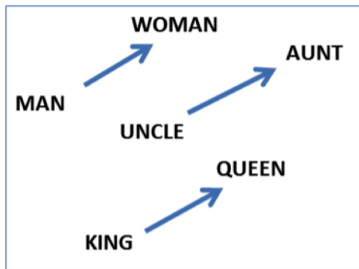
king = $\langle 55, -10, 176, 27 \rangle$

man = $\langle 10, 79, 150, 83 \rangle$

woman = $\langle 15, 74, 159, 106 \rangle$

queen = $\langle 60, -15, 185, 50 \rangle$

Word Embeddings



$$\text{vec}(\text{King}) - \text{vec}(\text{Man}) + \text{vec}(\text{Woman}) \approx \text{vec}(\text{Queen})^2$$

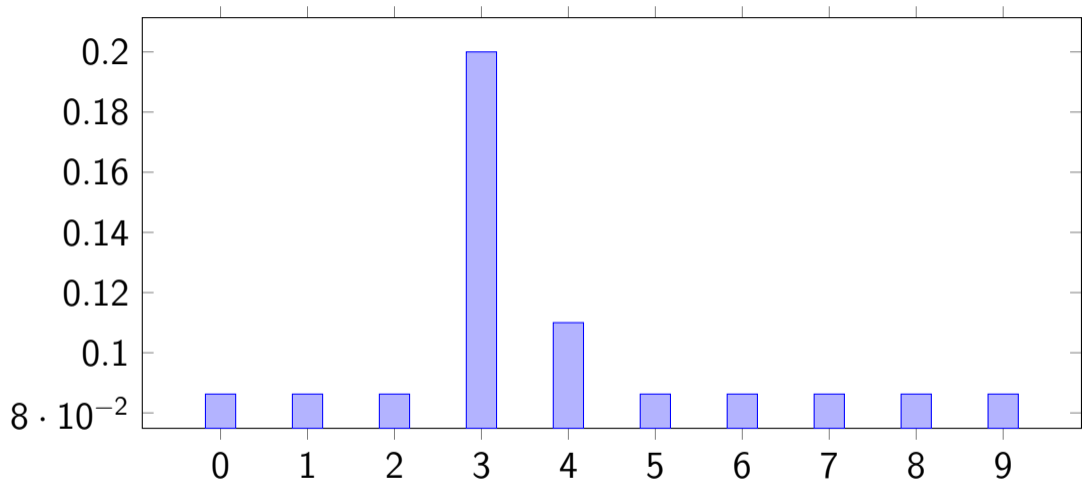
²Linguistic Regularities in Continuous Space Word Representations Mikolov et al. (2013)

Language Models

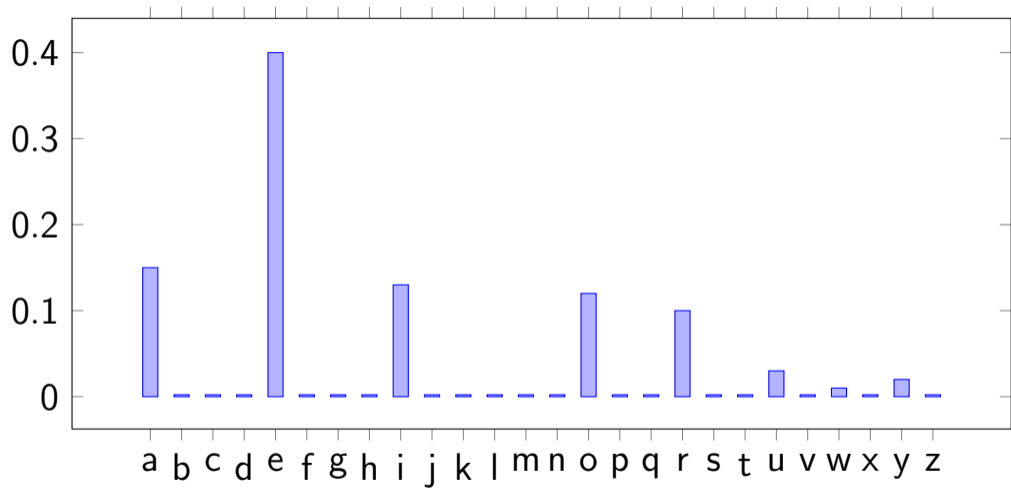
Language Models

- ▶ Language is **sequential** and has lots of words.

1,2,?



Th?



- ▶ A language model can compute:
 1. the probability of an upcoming symbol:

$$P(w_n | w_1, \dots, w_{n-1})$$

2. the probability for a sequence of symbols³

$$P(w_1, \dots, w_n)$$

³We can go from 1. to 2. using the Chain Rule of Probability $P(w_1, w_2, w_3) = P(w_1)P(w_2|w_1)P(w_3|w_1, w_2)$

- ▶ Language models are useful for machine translation because they help with:

1. word ordering

$$P(\textit{Yes I can help you}) > P(\textit{Help you I can yes})^4$$

2. word choice

$$P(\textit{Feel the Force}) > P(\textit{Eat the Force})$$

⁴Unless its Yoda that speaking

Neural Language Models

Recurrent Neural Networks

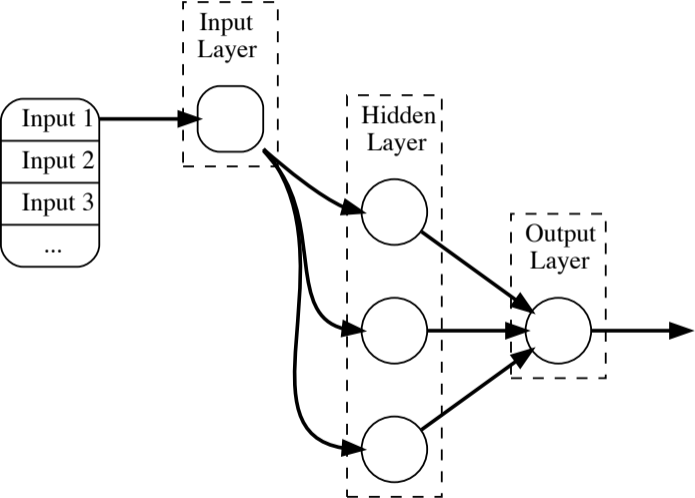
A particular type of neural network that is useful for processing **sequential** data (such as, language) is a **RECURRENT NEURAL NETWORK**.

Recurrent Neural Networks

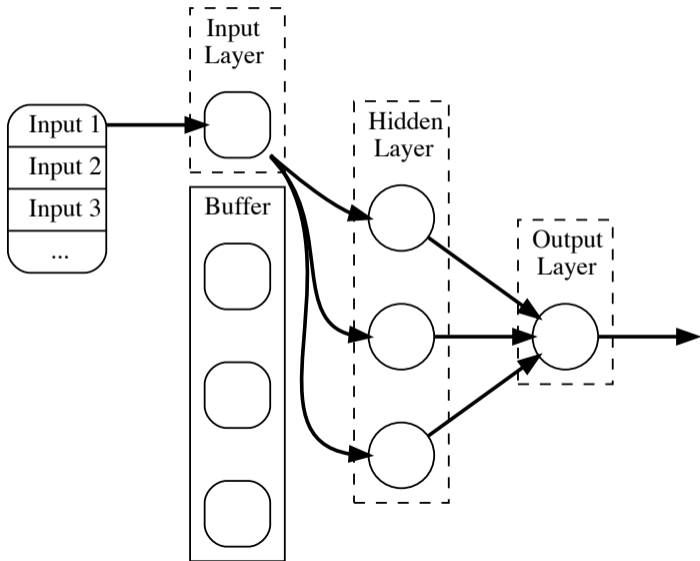
Using an RNN we process our sequential data one input at a time.

In an RNN the outputs of some of the neurons for one input are feed back into the network as part the next input.

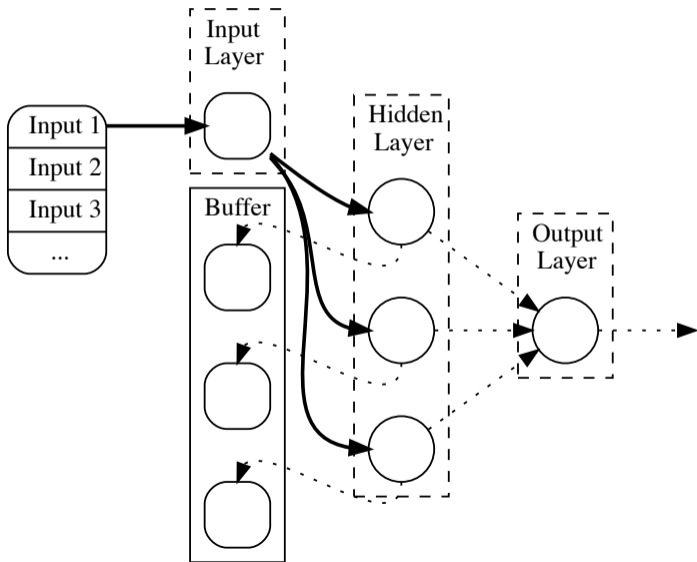
Simple Feed-Forward Network



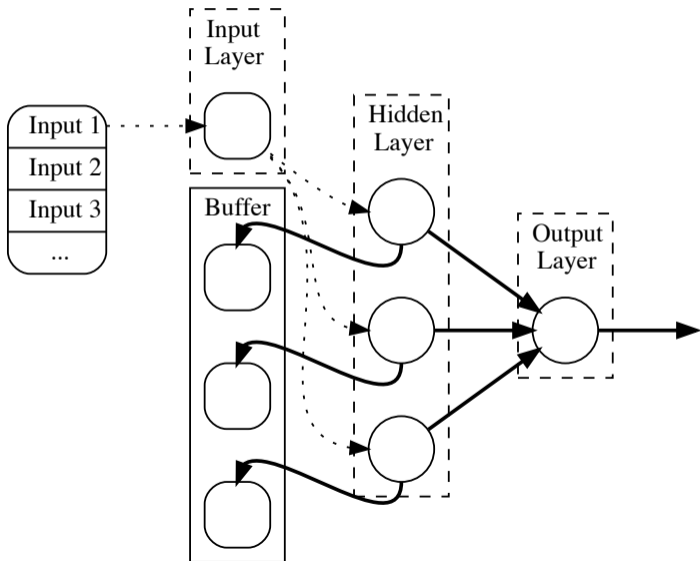
Recurrent Neural Networks



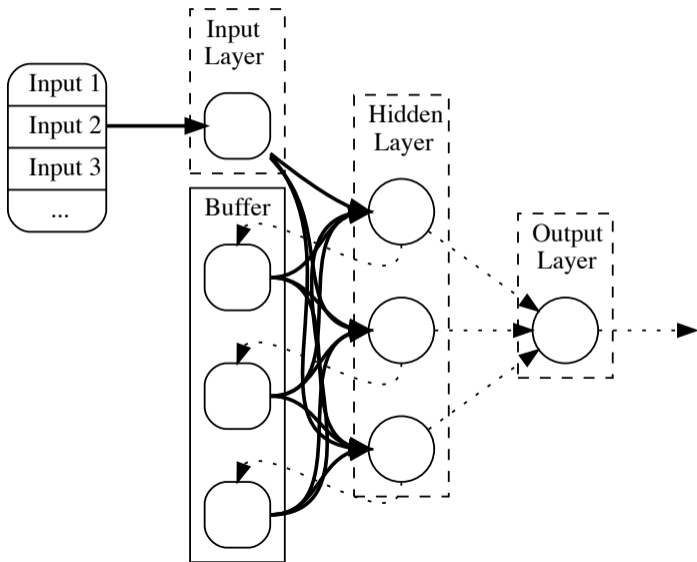
Recurrent Neural Networks



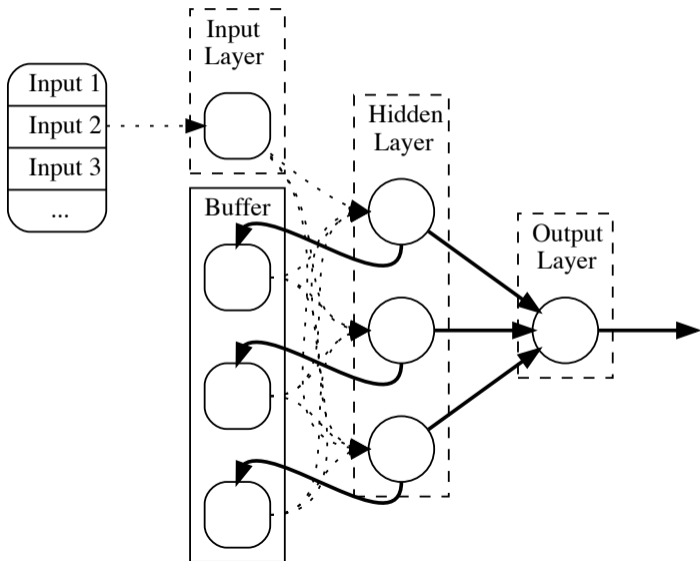
Recurrent Neural Networks



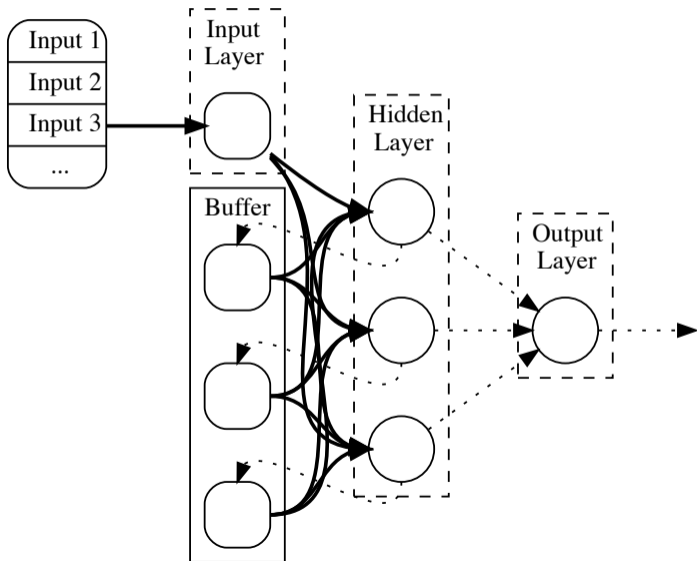
Recurrent Neural Networks



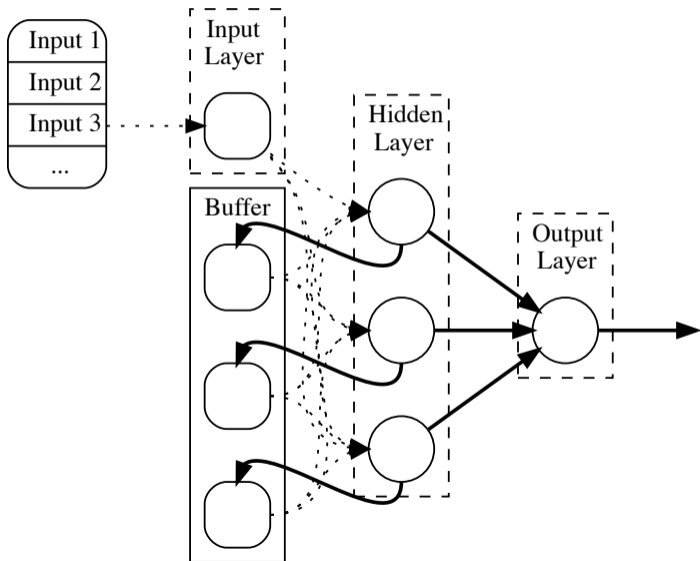
Recurrent Neural Networks

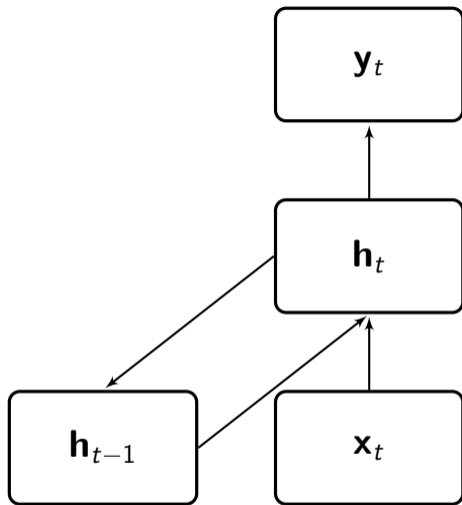


Recurrent Neural Networks



Recurrent Neural Networks





$$\mathbf{h}_t = \phi((\mathbf{W}_{hh} \cdot \mathbf{h}_{t-1}) + (\mathbf{W}_{xh} \cdot \mathbf{x}_t))$$

$$\mathbf{y}_t = \phi(\mathbf{W}_{hy} \cdot \mathbf{h}_t)$$

Figure: Recurrent Neural Network

Recurrent Neural Networks

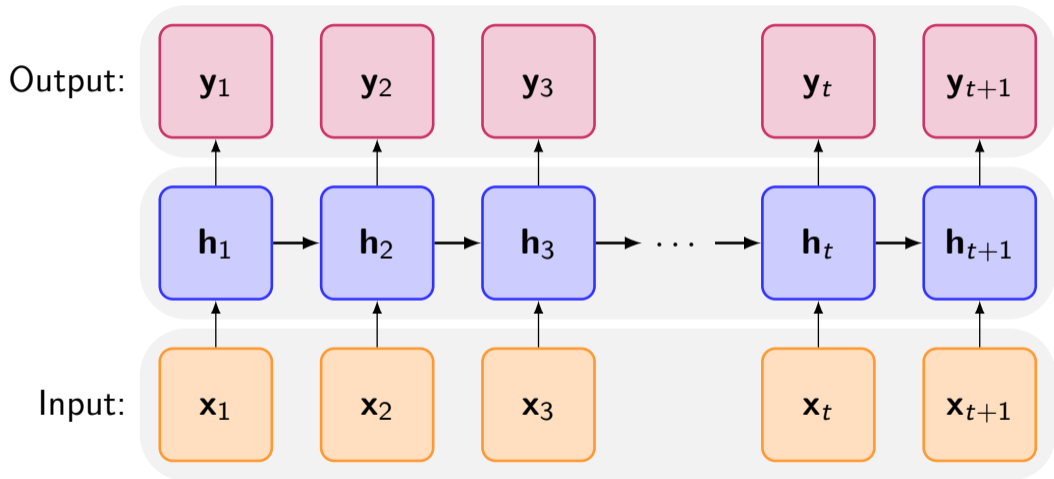
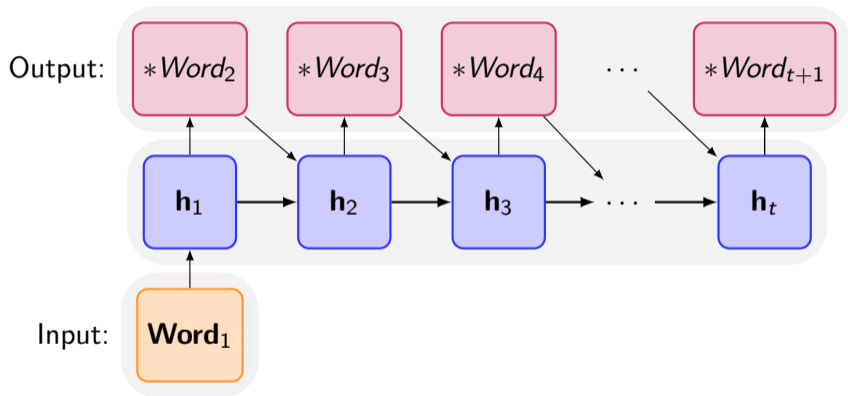


Figure: RNN Unrolled Through Time

Hallucinating Text



Hallucinating Shakespeare

PANDARUS: Alas, I think he shall be come approached and the day
When little strain would be attain'd into being never fed, And who is
but a chain and subjects of his death, I should not sleep.

Second Senator: They are away this miseries, produced upon my
soul, Breaking and strongly should be buried, when I perish The
earth and thoughts of many states.

DUKE VINCENTIO: Well, your wit is in the care of side and that.

From: <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

Neural Machine Translation

Neural Machine Translation

1. RNN Encoders
2. RNN Language Models

Encoders

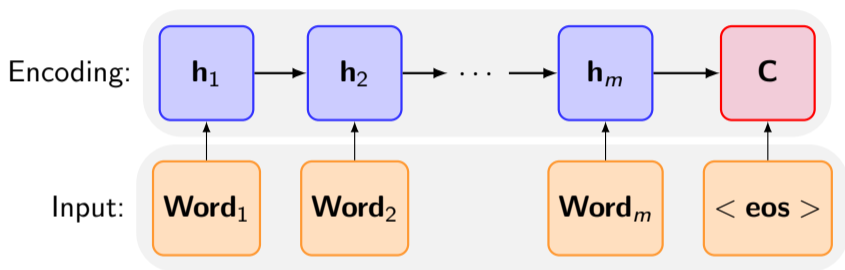


Figure: Using an RNN to Generate an Encoding of a Word Sequence

Language Models

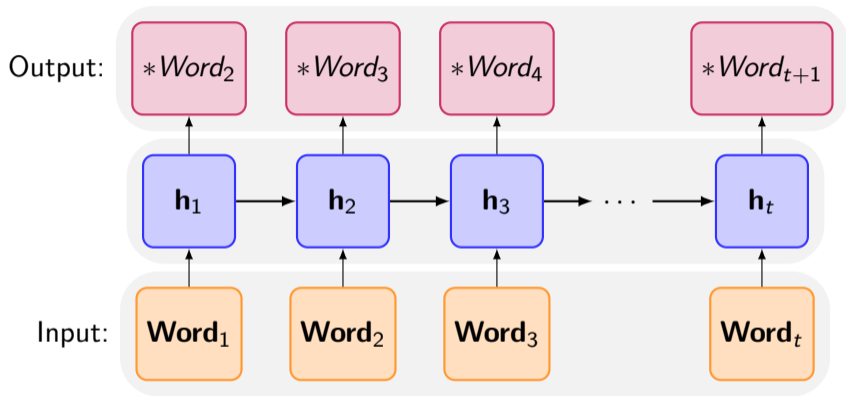


Figure: RNN Language Model Unrolled Through Time

Decoder

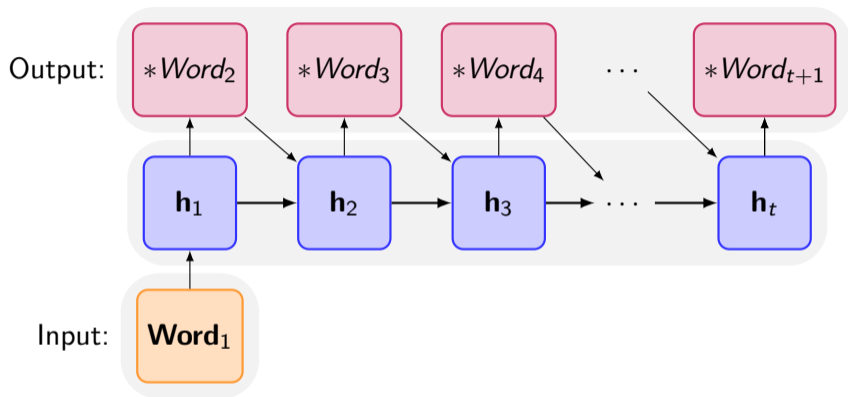


Figure: Using an RNN Language Model to Generate (Hallucinate) a Word Sequence

Encoder-Decoder Architecture

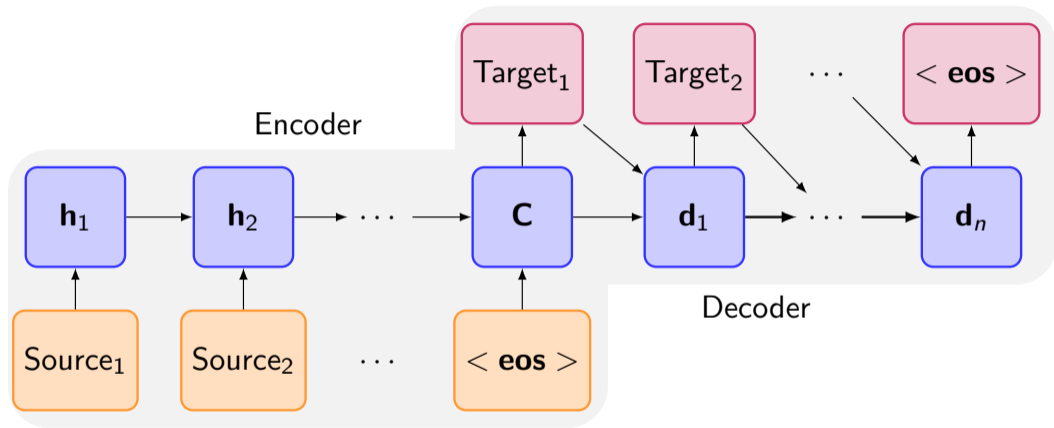


Figure: Sequence to Sequence Translation using an Encoder-Decoder Architecture

Neural Machine Translation

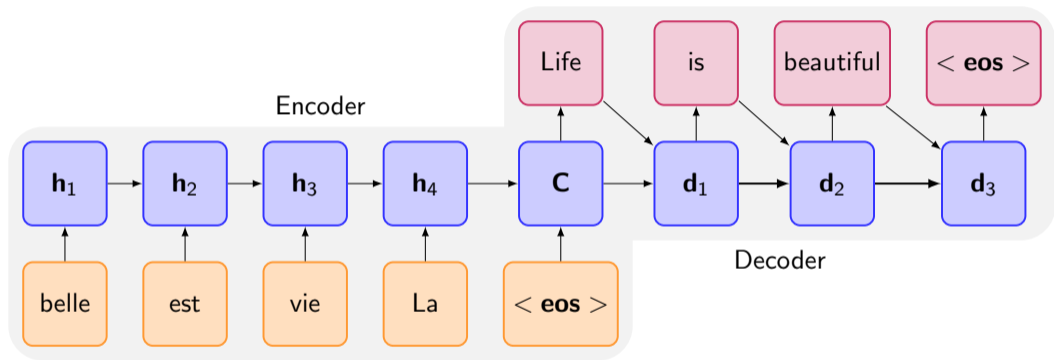


Figure: Example Translation using an Encoder-Decoder Architecture

Beyond NMT: Image Annotation

Image Annotation

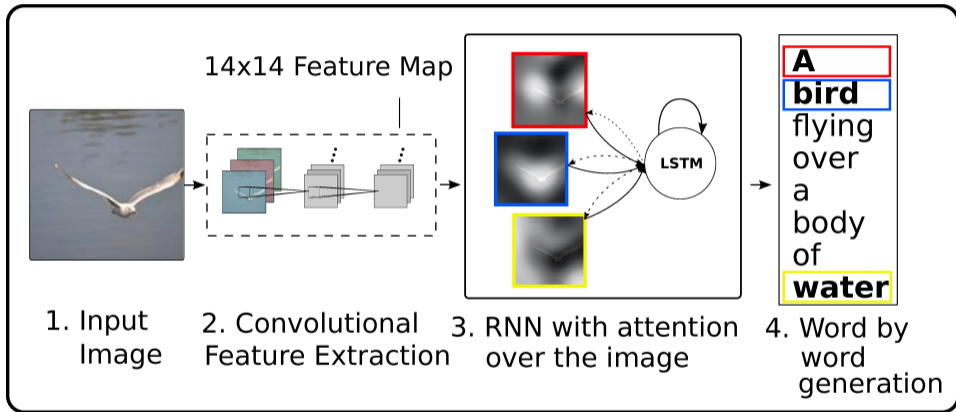


Image from Image from Show, Attend and Tell: Neural Image Caption Generation with Visual Attention Xu et al. (2015)

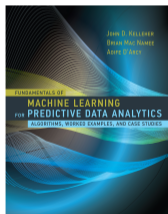
Thank you for your attention

john.d.kelleher@dit.ie

@johndkelleher

www.machinelearningbook.com

<https://ie.linkedin.com/in/johndkelleher>



Acknowledgements: The ADAPT Centre is funded under the SFI Research Centres Programme (Grant 13/RC/2106) and is co-funded under the European Regional Development Fund.

Books: Kelleher et al. (2015) Kelleher and Tierney (2018)

References I

- Kelleher, J. (2016). Fundamentals of machine learning for neural machine translation. In *Translating Europe Forum 2016: Focusing on Translation Technologies*. The European Commission Directorate-General for Translation, European Commission doi: 10.21427/D78012.
- Kelleher, J. D., Mac Namee, B., and D'Arcy, A. (2015). *Fundamentals of Machine Learning for Predictive Analytics: Algorithms, Worked Examples and Case Studies*. MIT Press, Cambridge, MA.
- Kelleher, J. D. and Tierney, B. (2018). *Data Science*. MIT Press.

References II

- Mikolov, T., Yih, W.-t., and Zweig, G. (2013). Linguistic regularities in continuous space word representations. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 746–751.
- Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhudinov, R., Zemel, R., and Bengio, Y. (2015). Show, attend and tell: Neural image caption generation with visual attention. In *International Conference on Machine Learning*, pages 2048–2057.