

---

# Advanced Deep Learning: Recurrent Networks

Anders Søgaard

---

coASfal

# Course outline

**Goal 1:** Quick tour of recent developments in deep learning

**Goal 2:** Inspiration for thesis/research projects

Week	Lecturer	Subject	Literature	Assignment
1	Stefan	Introduction to Neural Networks.	d2l 2.1-2.5, 2.7, 11.5.1, <b>slides</b>	
2	Stefan	CNNs; FCNs; U-Nets. Data augmentation; invariance; regularization e.g. dropout	d2l 6, 7, <b>slides</b>	Assignment 1 (May 10)
3	Anders/Phillip	May 9 (A): RNNs May 11 (P): Transformers	d2l 8 Transformers: d2l 10.5-10.7 + <a href="#">Vaswani et al. (2017)</a> <sup>↗</sup>	Assignment 2 (May 17)
4	Phillip/Anders	May 16 (P): Representation and Adversarial Learning May 18 (A): A Learning Framework + Self-supervised Learning + Contrastive Learning	Autoencoders: tba GANs: tba Self-supervised learning: <a href="#">blog post</a> <sup>↗</sup> Contrastive learning: <a href="#">Dor et al. (2018)</a> <sup>↗</sup>	
5	Anders	May 23-25: Applications of Representation, Adversarial and Contrastive Learning	GANs: <a href="#">Lample et al. (2018)</a> <sup>↗</sup> Autoencoders: <a href="#">Chandar et al. (2011)</a> <sup>↗</sup> Contrastive learning: <a href="#">Yu et al. (2018)</a> <sup>↗</sup> Examples: <a href="#">Goodfellow et al. (2015)</a> <sup>↗</sup> DynaBench: <a href="#">Talk by Douwe Kiela</a> <sup>↗</sup> (Facebook, now HuggingFace)	Assignment 3 [MC on Representation Learning/1p Report on Adversarial Learning] (May 31)
6	Anders	May 30-June 1: Interpretability	Literature: <a href="#">Sogaard (2022)</a> <sup>↗</sup>	
7	Anders	June 8: Interpretability (Note: June 6 off)	Literature: <a href="#">Feng and Boyd-Graber (2018)</a> <sup>↗</sup> ; <a href="#">Jiang and Senge (2021)</a> <sup>↗</sup>	
8	Anders	June 13-15: Best Practices	Literature: <a href="#">Dodge et al. (2019)</a> <sup>↗</sup> and <a href="#">Raji et al. (2021)</a> <sup>↗</sup>	Assignment 4 [MC on Interpretability; 1p Report on Best Practices] (June 21)

# Course outline

**Goal 1:** Quick tour of recent developments in deep learning

**Goal 2:** Inspiration for thesis/research projects

**Stefan:** U-Nets+CNNs (1-2)

**Phillip and I:** a) RNNs and Transformers and how to train them (3-5); b) Interpretability and Best Practices (6-8)

Week	Lecturer	Subject	Literature	Assignment
1	Stefan	Introduction to Neural Networks.	d2l 2.1-2.5, 2.7, 11.5.1, <b>slides</b>	
2	Stefan	CNNs; FCNs; U-Nets. Data augmentation; invariance; regularization e.g. dropout	d2l 6, 7, <b>slides</b>	Assignment 1 (May 10)
3	Anders/Phillip	May 9 (A): RNNs May 11 (P): Transformers	d2l 8 Transformers: d2l 10.5-10.7 + <a href="#">Vaswani et al. (2017)</a> <sup>↗</sup>	Assignment 2 (May 17)
4	Phillip/Anders	May 16 (P): Representation and Adversarial Learning May 18 (A): A Learning Framework + Self-supervised Learning + Contrastive Learning	Autoencoders: tba GANs: tba Self-supervised learning: <a href="#">blog post</a> <sup>↗</sup> Contrastive learning: <a href="#">Dor et al. (2018)</a> <sup>↗</sup>	
5	Anders	May 23-25: Applications of Representation, Adversarial and Contrastive Learning	GANs: <a href="#">Lample et al. (2018)</a> <sup>↗</sup> Autoencoders: <a href="#">Chandar et al. (2011)</a> <sup>↗</sup> Contrastive learning: <a href="#">Yu et al. (2018)</a> <sup>↗</sup> Examples: <a href="#">Goodfellow et al. (2015)</a> <sup>↗</sup> DynaBench: <a href="#">Talk by Douwe Kiela</a> <sup>↗</sup> (Facebook, now HuggingFace)	Assignment 3 [ <i>MC on Representation Learning/1p Report on Adversarial Learning</i> ] (May 31)
6	Anders	May 30-June 1: Interpretability	Literature: <a href="#">Søgaard (2022)</a> <sup>↗</sup>	
7	Anders	June 8: Interpretability (Note: June 6 off)	Literature: <a href="#">Feng and Boyd-Graber (2018)</a> <sup>↗</sup> ; <a href="#">Jiang and Senge (2021)</a> <sup>↗</sup>	
8	Anders	June 13-15: Best Practices	Literature: <a href="#">Dodge et al. (2019)</a> <sup>↗</sup> and <a href="#">Raji et al. (2021)</a> <sup>↗</sup>	Assignment 4 [ <i>MC on Interpretability; 1p Report on Best Practices</i> ] (June 21)

---

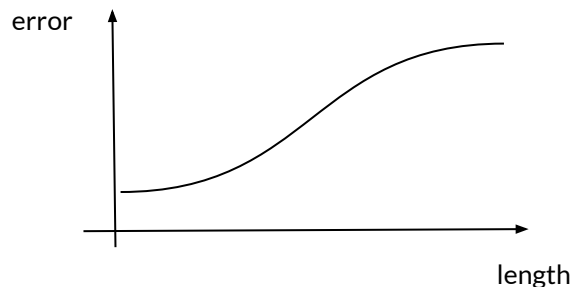
## Exercise (May 17)

- a) Multiple choice (Google Form)
  - b) 1 page report: Compare RNNs and LSTMs on  $a^n b^n$  and  $a^n b^n c^n$
  - c) ~~Code piece: RNN implementation~~
-

---

# Exercise (May 17)

- a) Multiple choice (Google Form)
- b) 1 page report: Compare RNNs and LSTMs on  $a^n b^n$  and  $a^n b^n c^n$
- c) ~~Code piece: RNN implementation~~



---

---

# Today

- a) The Deep Learning Landscape
  - b) Applications
  - c) Recurrency
  - d) Historical context
-

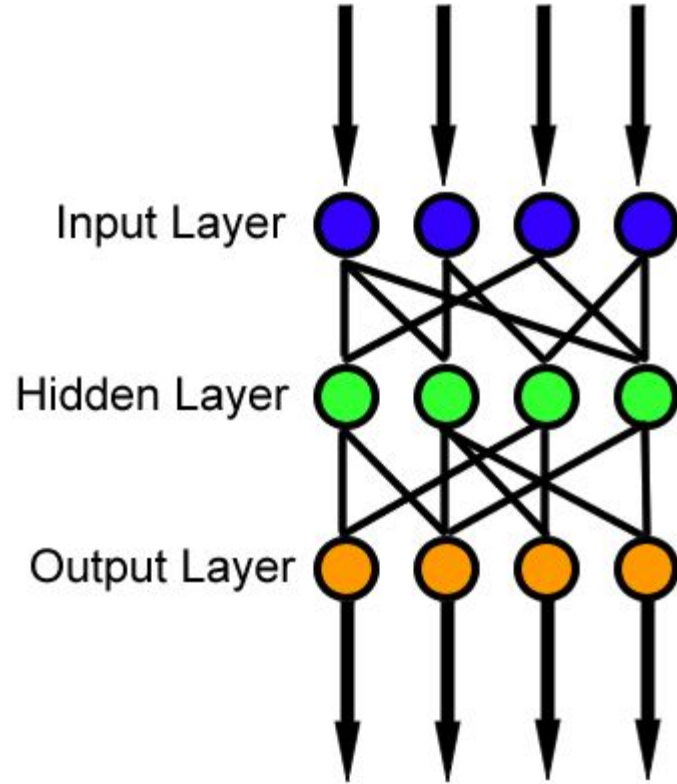
---

# Landscape

---

# Feed-forward Networks

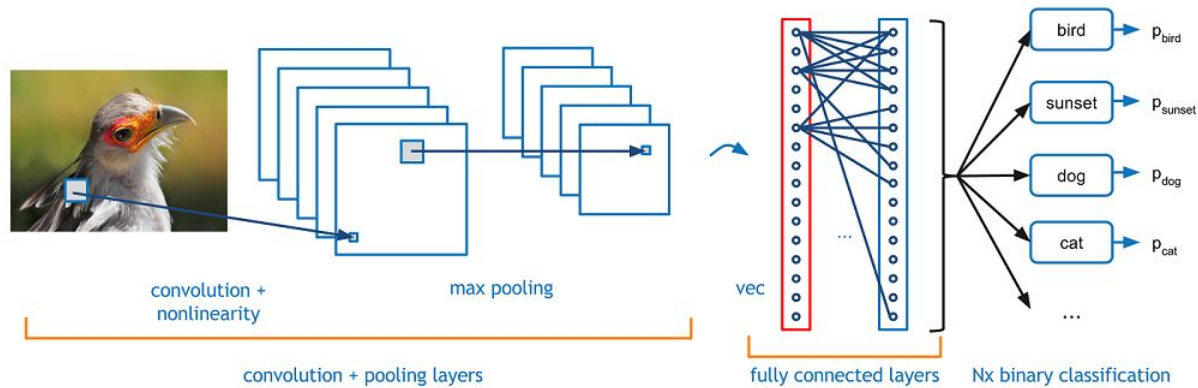
FFNs are networks of perceptrons or logistic regressors feeding into other perceptrons or logistic regressors. In the so-called *forward* step, that's the whole story. For training, you need to do a more complicated *backward* step to compute the appropriate weight updates. This process is called *back-propagation*.





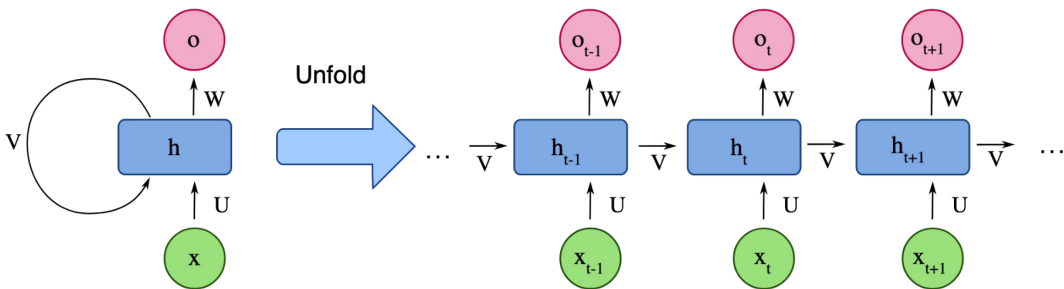
# Convolutional Neural Networks

Adds special early layers to FNNs that account for the invariance properties of images.



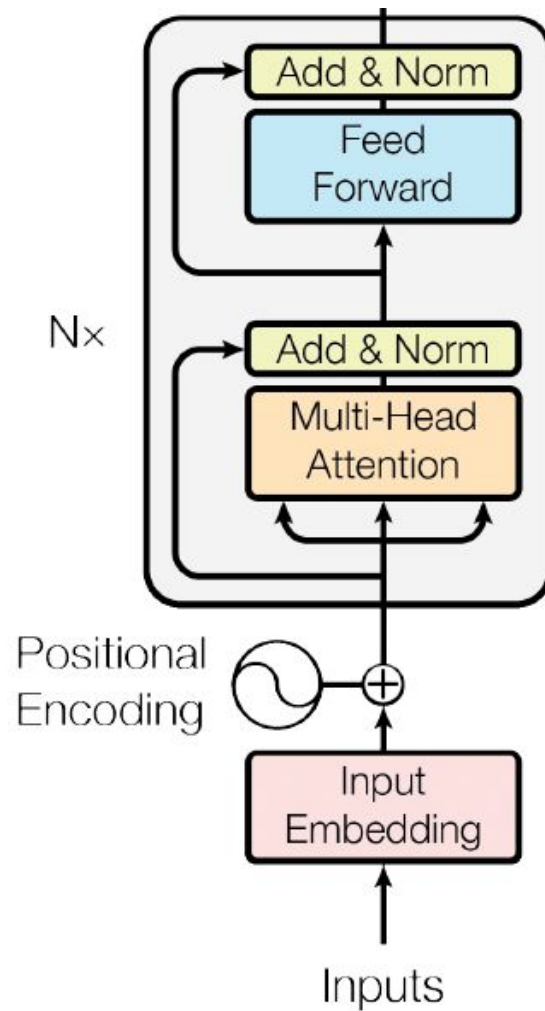
# Recurrent Neural Networks

You can think of an RNN as a FNN with  $n$  layers used to process a sequence of  $n$  tokens, but in which the  $n$  layers are all the same. **Note:** For a deep  $k$ -layer RNN, this would be  $kn$  layers with the  $n$  combinations of  $k$  layers being the same.



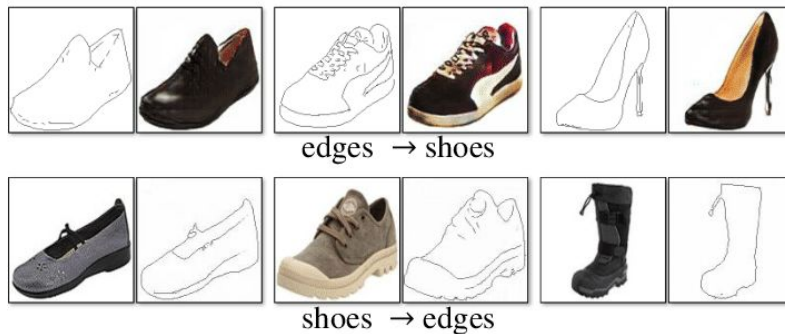
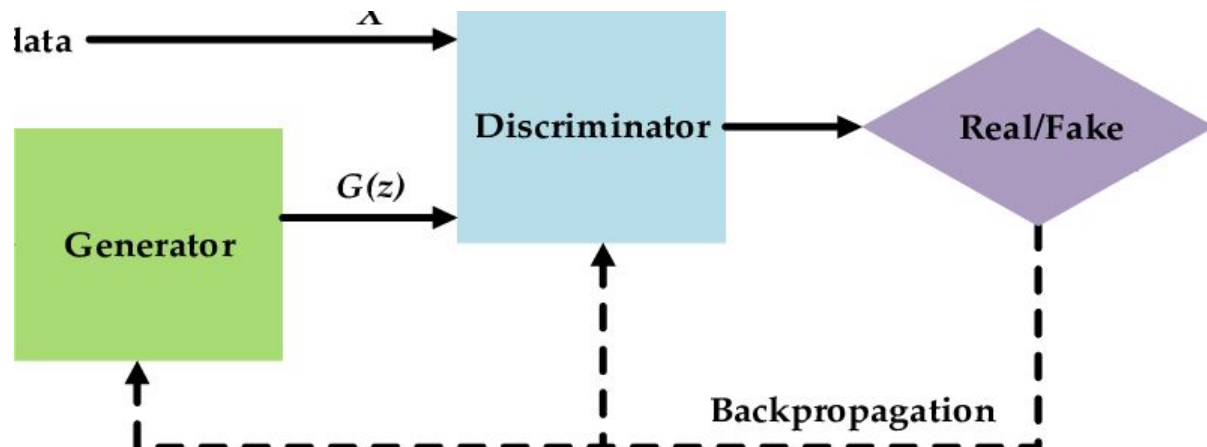
# Transformers

Improvement over RNNs, specialized for GPUs. Processes tokens in parallel, but considers their interaction with all other tokens (making inference quadratic-time).



# Generative Adversarial Networks

Combines two networks, e.g., a FNN and an RNN to transform data points from one distribution so they look like data points from another distribution.



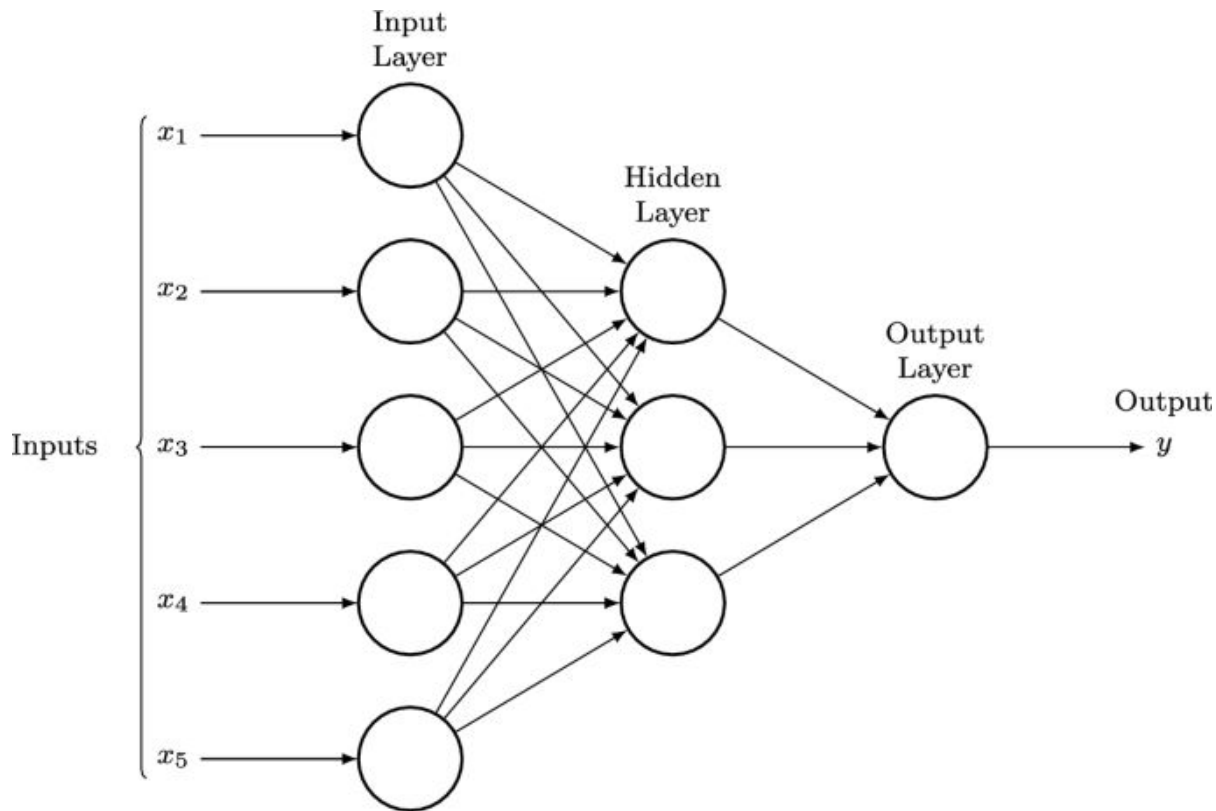
---

# Applications

---

# Sequential data

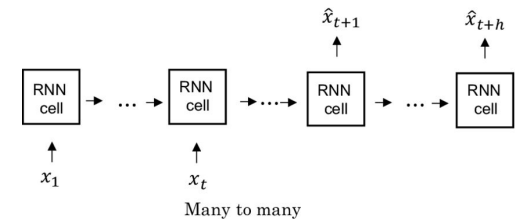
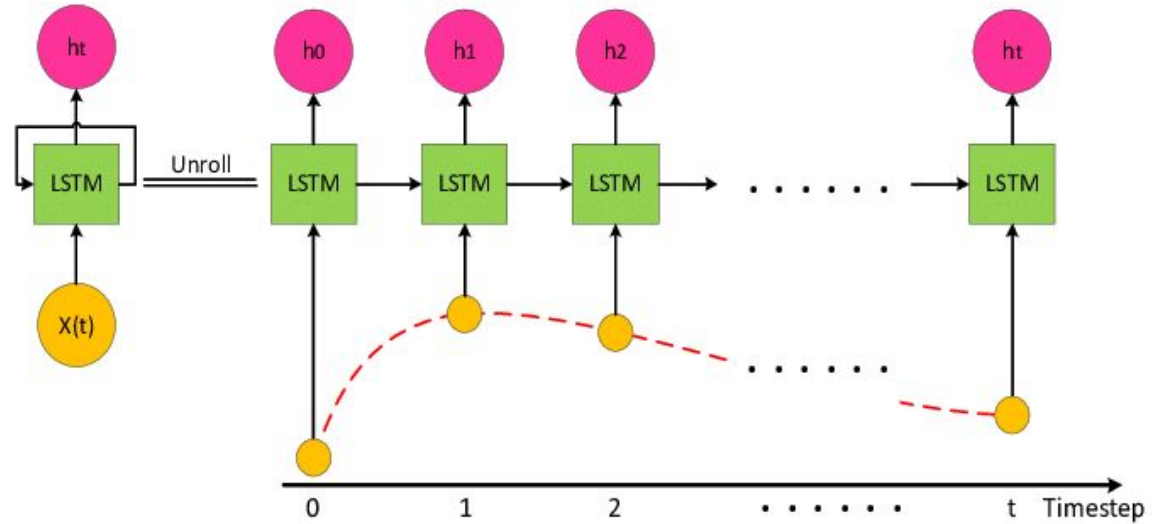
- Bag-of- $n$ -grams with feed-forward networks are limited by a poor bias-variance trade-off.
- So are 1D CNNs.
- Pre-deep models with slightly better trade-offs: HMMs, CRFs, and [HMM-Perceptrons](#).



# Time series data

Predicting trends or developments over time.

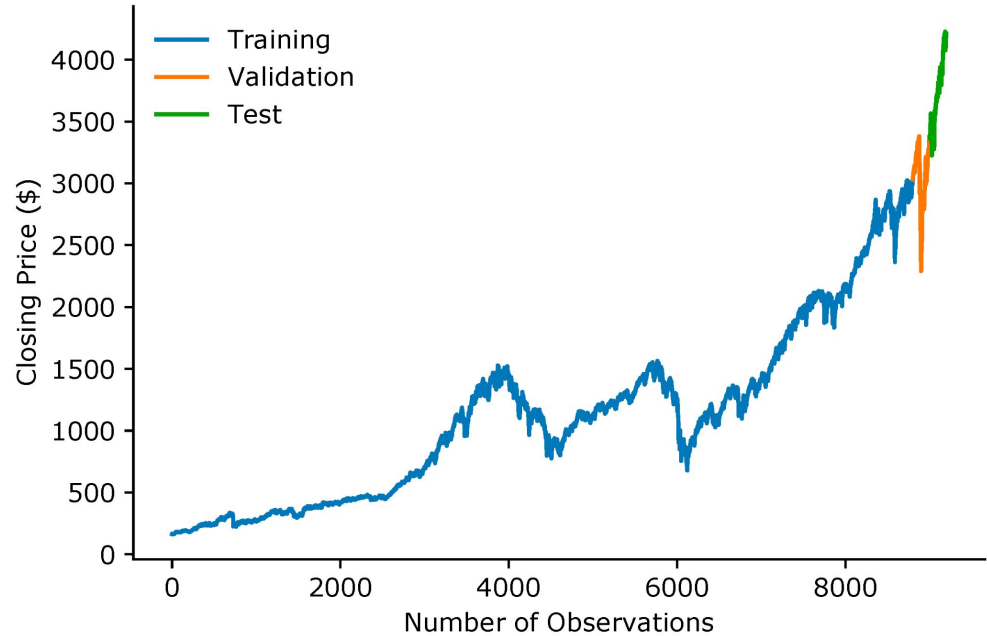
- a) Sales
- b) Climate
- c) Socio-economic variables
- d) ...



# Time series data

Predicting trends or developments over time.

- a) Sales
- b) Climate
- c) Socio-economic variables
- d) ...

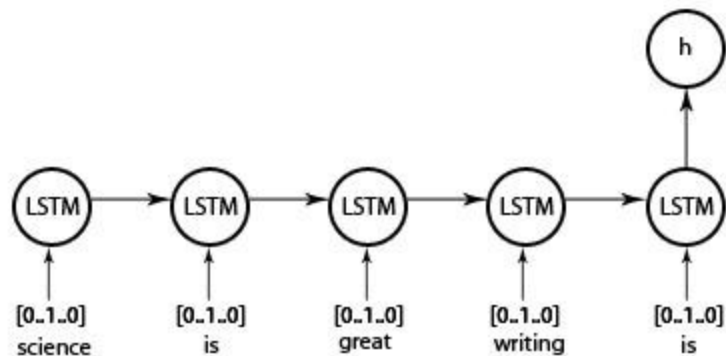


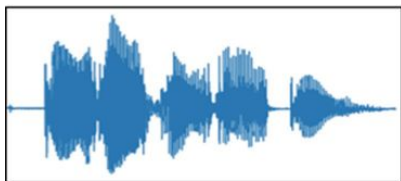


## Sentence or document classification

Sentences and documents come in different lengths.

- a) Sparse, e.g., few sentences of length 54
- b) Unbounded, e.g., may always potentially see a longer sentence





Raw Speech Signal

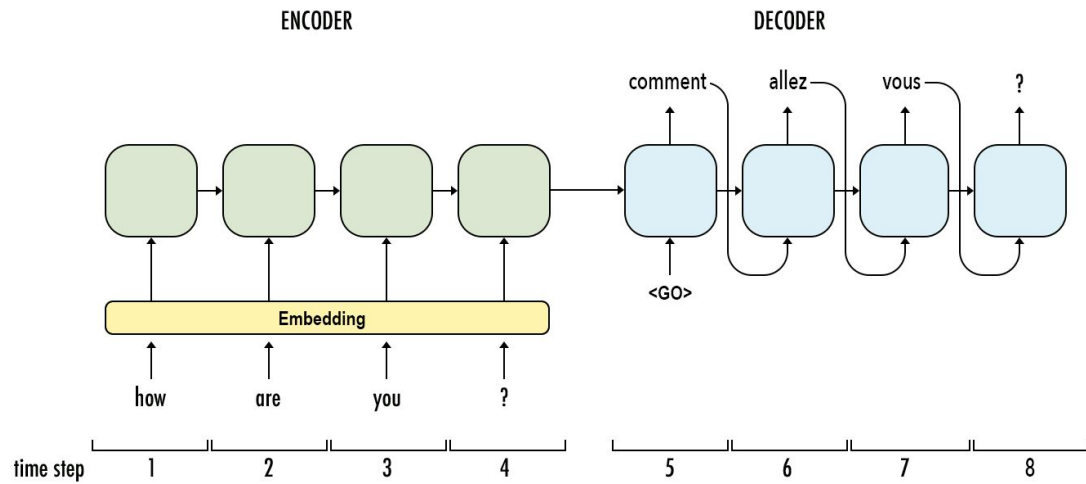


Do	you	understand	me
----	-----	------------	----

Transcription

Speech recognition

---



## Machine translation

---

# Examples of sequence data

Speech recognition



“The quick brown fox jumped  
over the lazy dog.”

Music generation



Sentiment classification

“There is nothing to like  
in this movie.”



DNA sequence analysis

AGCCCCTGTGAGGAACTAG



AGCCCCTGTGAGGAACTAG

Machine translation

Voulez-vous chanter avec  
moi?



Do you want to sing with  
me?

Video activity recognition



Running

Name entity recognition

Yesterday, Harry Potter  
met Hermione Granger.



Yesterday, Harry Potter  
met Hermione Granger.  
Andrew Ng

## Other examples

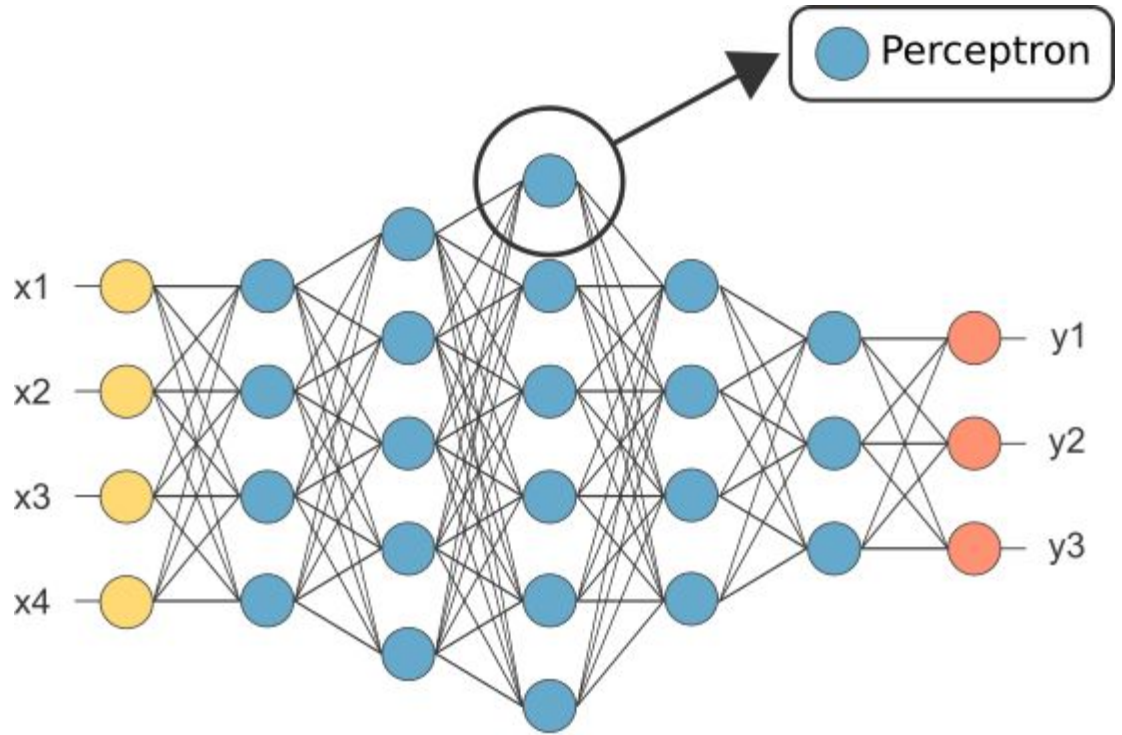
---

# Recurrency

---

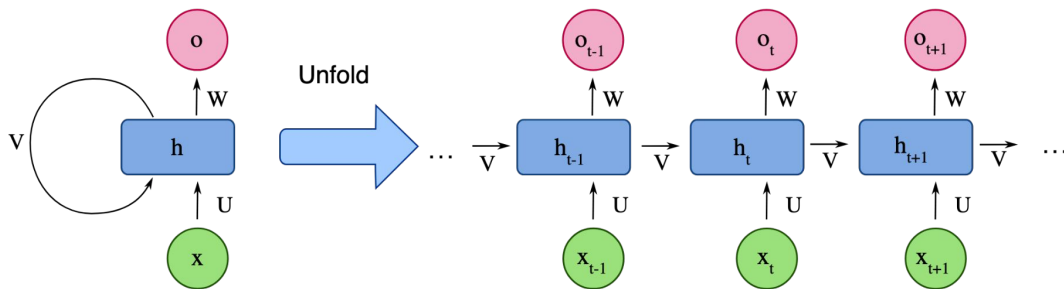
# Feed Forward

- a) Interconnected neurons
- b)  $m$  many layers
- c) But can't model  $n$  length sequences in an unbiased way



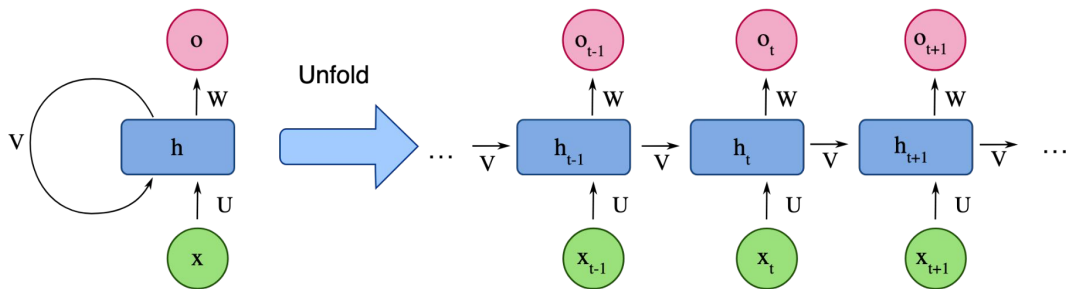
# Recurrent

But what if we used some of the hidden layers *recurrently*?



# Recurrent

But what if we used some of the hidden layers *recurrently*?



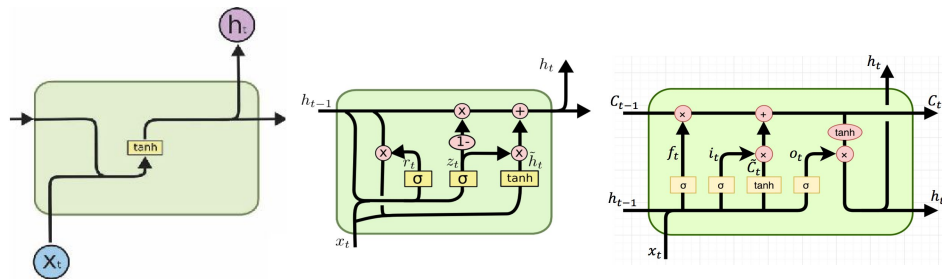
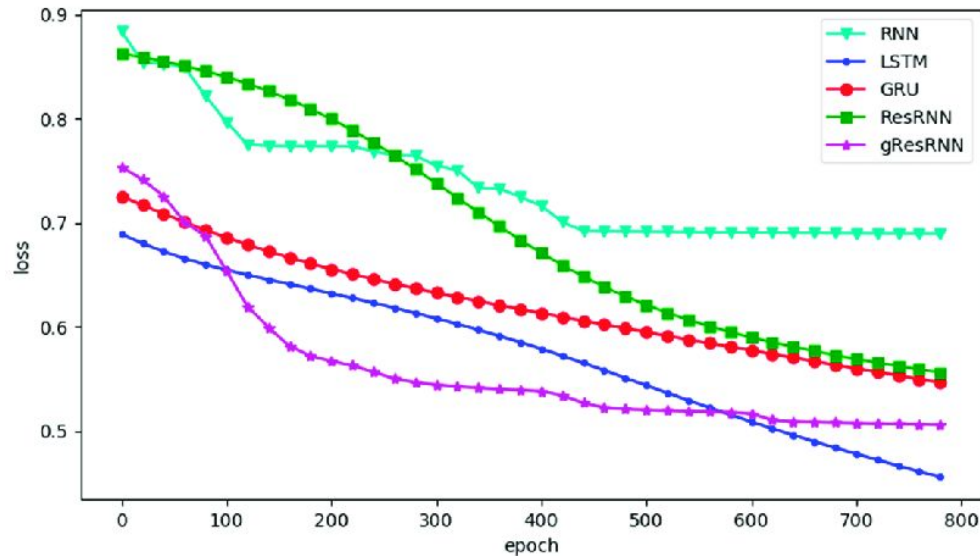
$$\begin{aligned} \mathbf{a}^{(t)} &= \mathbf{b} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)} \\ \mathbf{h}^{(t)} &= \tanh(\mathbf{a}^{(t)}) \\ \mathbf{o}^{(t)} &= \mathbf{c} + \mathbf{V}\mathbf{h}^{(t)} \\ \hat{\mathbf{y}}^{(t)} &= \text{softmax}(\mathbf{o}^{(t)}) \end{aligned}$$



# Beyond RNNs

- a) RNNs are basically stacked FNNs with parameter sharing.
- b) GRUs introduce update gates.
- c) LSTMs also introduce forget and output gates.

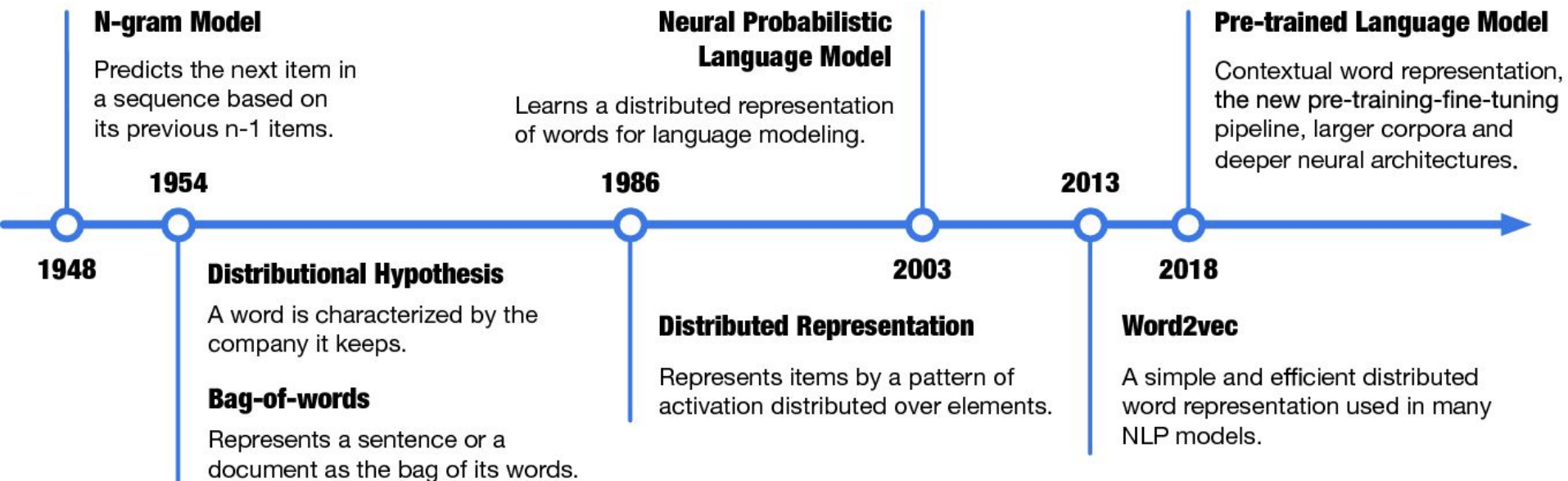
More control prevents vanishing gradients and enables more efficient induction of long-distance dependencies.



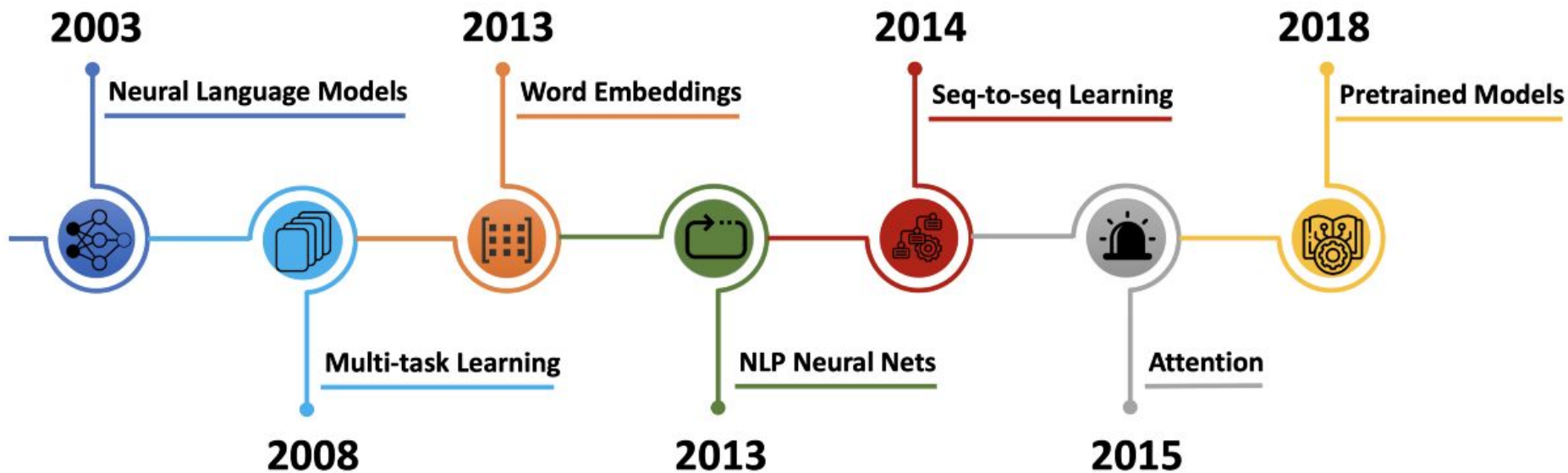
---

# Historical context

---



## History of NLP

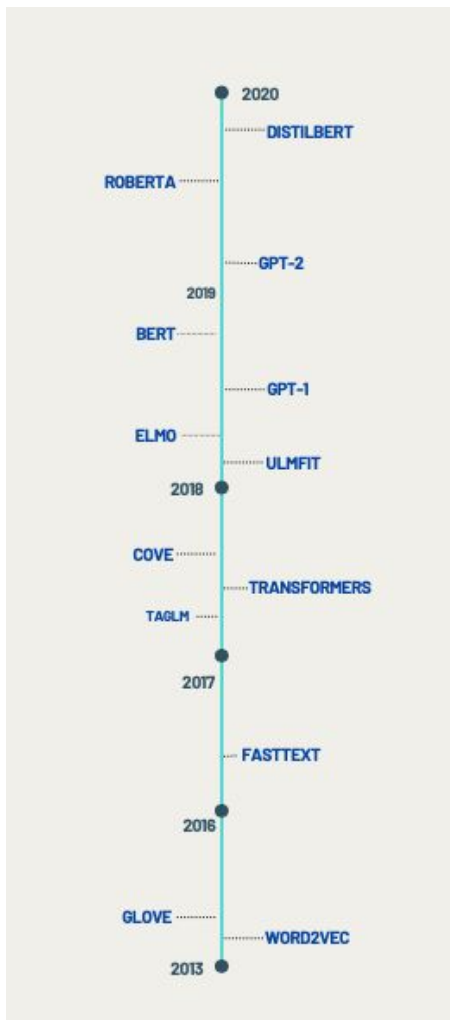


## History of NLP

---

# Pretrained models

- a) Early generations (w2v, fastText, etc) built on feed forward networks
- b) Intermediate models were based on RNNs (ELMO, GPT-1, etc)
- c) Later generations (GPT-3, BART, Pegasus, RoBERTa, etc) are based on Transformers (**next time**)



—

	<b>RNNs</b>	<b>GRUs/LSTMs</b>	<b>Transformers</b>
<b>NLP</b>	2010	2013	2017
<b>Computer Vision</b>	(Video)	(Video)	ViT (2021) TrOCR (2021)

## History of NLP and Computer Vision

---