

# Deep Learning for NLP

- Oliver Guhr  
I833 SS2019

# Natural Language Processing

# NLP Tasks

- Easy
  - Spell Checking
  - Keyword Search

# NLP Tasks

- Easy
  - Spell Checking
  - Keyword Search
- Medium
  - Parsing information from unstructured text

# NLP Tasks

- Easy
  - Spell Checking
  - Keyword Search
- Medium
  - Parsing information from unstructured text
- Hard
  - Machine Translation
  - Semantic Analysis
  - Question Answering

Übersetze **Englisch** (erkannt)

A very simple way to improve the performance of almost any machine learning algorithm is to train many different models on the same data and then to average their predictions. Unfortunately, making predictions using a whole ensemble of models is cumbersome and may be too computationally expensive to allow deployment to a large number of users, especially if the individual models are large neural nets. Caruana and his collaborators have shown that it is possible to compress the knowledge in an ensemble into a single model which is much easier to deploy and we develop this approach further using a different compression technique. We achieve some surprising results on MNIST and we show that we can significantly improve the acoustic model of a heavily used commercial system by distilling the knowledge in an ensemble of models into a single model. We also introduce a new type of ensemble composed of one or more full models and many specialist models which learn to distinguish fine-grained classes that the full models confuse. Unlike a mixture of experts, these specialist models can be trained rapidly and in parallel.

Dokument übersetzen

1130/5000

Übersetze nach **Deutsch**

Eine sehr einfache Möglichkeit, die Leistung fast jedes maschinellen Lernalgorithmus zu verbessern, besteht darin, viele verschiedene Modelle auf den gleichen Daten zu trainieren und dann ihre Vorhersagen zu berechnen. Leider ist die Vorhersage mit einem ganzen Ensemble von Modellen umständlich und kann zu rechenintensiv sein, um die Bereitstellung für eine große Anzahl von Benutzern zu ermöglichen, insbesondere wenn die einzelnen Modelle große neuronale Netze sind. Caruana und seine Mitarbeiter haben gezeigt, dass es möglich ist, das Wissen in einem Ensemble zu einem einzigen, viel einfacher zu implementierenden Modell zu komprimieren, und wir entwickeln diesen Ansatz mit einer anderen Kompressionstechnik weiter. Wir erzielen einige überraschende Ergebnisse auf MNIST und zeigen, dass wir das akustische Modell eines stark genutzten kommerziellen Systems signifikant verbessern können, indem wir das Wissen in einem Ensemble von Modellen in ein einziges Modell überführen. Wir stellen auch eine neue Art von Ensemble vor, das aus einem oder mehreren Vollmodellen und vielen Spezialmodellen besteht, die lernen, feinkörnige Klassen zu unterscheiden, die die Vollmodelle verwirren. Im Gegensatz zu einer Mischung von Experten können diese Spezialmodelle schnell und parallel trainiert werden.

**Experte** noun, masculine

**expert** n (plural: **experts**)

Experten bewerteten die Qualität der Produkte.

Die Maschinen werden von Experten kontrolliert.

Ein Gremium aus Experten berät die Regierung.

Experts evaluated the quality of the products.

The machines will be inspected by experts.

A body of experts is advising the government.

**specialist** n (plural: **specialists**)



Artificial intelligence (AI)

Alex Hern @alexhern Thu 14 Feb 2019 17.00 GMT

6.473 572

New AI fake text generator may be too dangerous to release, say creators

The Elon Musk-backed nonprofit company OpenAI declines to release research publicly for fear of misuse



Editorially independent, open to everyone

We chose a different approach – will you support it?

Support The Guardian →

most viewed

- Live US-China trade war: Beijing vows to retaliate as tariffs raised - Business live
- Anna Sorokin: fake German heiress sentenced to up to 12 years in prison
- Freddie Starr: comedian found dead at home in

Support

Available for every

Contribute →

News

World UK Science

Artificial intel  
(AI)

Login

INDY/LIFE

OpenAI built a text generator so good, it's  
considered too dangerous to release

Zack Whittaker @zackwhittaker / 3 months ago

OpenAI let us try its state-of-the-art NLP text  
generator

AI TEXT GENERATOR TOO  
DANGEROUS TO RELEASE, SAY  
CREATORS

Developers cite concerns over fake news proliferation and risk of online impersonation

ional edition ▾



ent,  
everyone

approach –



1 - Business

in: fake German  
enced to up to  
rison

r: comedian  
at home in



# The GTP-2 Model

- Model trained on 40GB text
- transformer-based language model
- Objective: predict the next word, given all of the previous words within some text
- Try (a smaller) version online:
  - <https://talktotransformer.com>
- Paper and Code:
  - <https://openai.com/blog/better-language-models/>

## **human-written input:**

*In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.*

## **model output:**

The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

[...]

# SQuAD2.0

The Stanford Question Answering Dataset

# 1973\_oil\_crisis

## The Stanford Question Answering Dataset

The 1973 oil crisis began in October 1973 when the members of the Organization of Arab Petroleum Exporting Countries (OAPEC, consisting of the Arab members of OPEC plus Egypt and Syria) proclaimed an oil embargo. By the end of the embargo in March 1974, the price of oil had risen from US\$3 per barrel to nearly \$12 globally; US prices were significantly higher. The embargo caused an oil crisis, or "shock", with many short- and long-term effects on global politics and the global economy. It was later called the "first oil shock", followed by the 1979 oil crisis, termed the "second oil shock."

**When did the 1973 oil crisis begin?**

*Ground Truth Answers:* October 1973 October 1973 October 1973 October 1973

**What was the price of oil in March of 1974?**

*Ground Truth Answers:* nearly \$12 \$12 \$12 \$12 \$12

**When was the second oil crisis?**

*Ground Truth Answers:* 1979 1979 1979 1979 1979

**What was another term used for the oil crisis?**

*Ground Truth Answers:* first oil shock shock shock first oil shock shock

**Who proclaimed the oil embargo?**

*Ground Truth Answers:* members of the Organization of Arab Petroleum Exporting Countries members of the Organization of Arab Petroleum Exporting Countries Organization of Arab Petroleum Exporting Countries members of the Organization of Arab Petroleum Exporting Countries OAPEC



# Leaderboard

Rank	Model	EM	F1
	Human Performance <i>Stanford University</i> ( <a href="#">Rajpurkar &amp; Jia et al. '18</a> )	86.831	89.452
1 Mar 20, 2019	BERT + DAE + AoA (ensemble) <i>Joint Laboratory of HIT and iFLYTEK Research</i>	87.147	89.474
2 Mar 15, 2019	BERT + ConvLSTM + MTL + Verifier (ensemble) <i>Layer 6 AI</i>	86.730	89.286
3 Mar 05, 2019	BERT + N-Gram Masking + Synthetic Self- Training (ensemble) <i>Google AI Language</i> <a href="https://github.com/google-research/bert">https://github.com/google-research/bert</a>	86.673	89.147

## Leaderboard

Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar & Jia et al. '18)	86.831	89.452
1 Mar 20, 2019	BERT + DAE + AoA (ensemble) Joint Laboratory of HIT and iFLYTEK Research	87.147	89.474
2 Mar 15, 2019	BERT + ConvLSTM + MTL + Verifier (ensemble) Layer 6 AI	86.730	89.286
3 Mar 05, 2019	BERT + N-Gram Masking + Synthetic Self-Training (ensemble) Google AI Language <a href="https://github.com/google-research/bert">https://github.com/google-research/bert</a>	86.673	89.147
4 Apr 13, 2019	SemBERT (ensemble) Shanghai Jiao Tong University	86.166	88.886
5 Mar 16, 2019	BERT + DAE + AoA (single model) Joint Laboratory of HIT and iFLYTEK Research	85.884	88.621
6 Mar 05, 2019	BERT + N-Gram Masking + Synthetic Self-Training (single model) Google AI Language <a href="https://github.com/google-research/bert">https://github.com/google-research/bert</a>	85.150	87.715
7 Jan 15, 2019	BERT + MMFT + ADA (ensemble) Microsoft Research Asia	85.082	87.615
7 Mar 13, 2019	BERT + ConvLSTM + MTL + Verifier (single model) Layer 6 AI	84.924	88.204
8 Apr 16, 2019	Insight-baseline-BERT (single model) PAI Insight Team	84.834	87.644

Rank	Task	Model	DP1	F1
		<b>Human Performance</b> Wikipedia (Human) / Wikipedia (AI) / 100	80.00	80.00
1	<b>GPT-4o</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	87.87	88.49
2		<b>Claude 3.5 Sonnet</b> (unlimited)	86.70	88.36
3	<b>Gemini 2.0 Flash</b> (unlimited)	Google AI / Gemini 2.0 Flash / 90B	86.50	87.87
4		<b>GPT-4o</b> (unlimited)	86.06	86.65
5	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	85.98	86.43
6		<b>GPT-4o</b> (unlimited)	85.90	87.78
7	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	85.00	87.00
8		<b>Claude 3.5 Sonnet</b> (unlimited)	84.58	86.36
9	<b>GPT-4o</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	84.50	87.60
10		<b>GPT-4o</b> (unlimited)	84.00	86.22
11	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.96	86.22
11		<b>Claude 3.5 Sonnet</b> (unlimited)	83.49	86.22
11	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.40	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12		<b>Claude 3.5 Sonnet</b> (unlimited)	83.36	86.22
12	<b>Claude 3.5 Sonnet</b> (unlimited)	OpenAI LLaMA v3.1 70B / 128B / 320B	83.36	86.22
12				

BERT will be  
part of the  
course



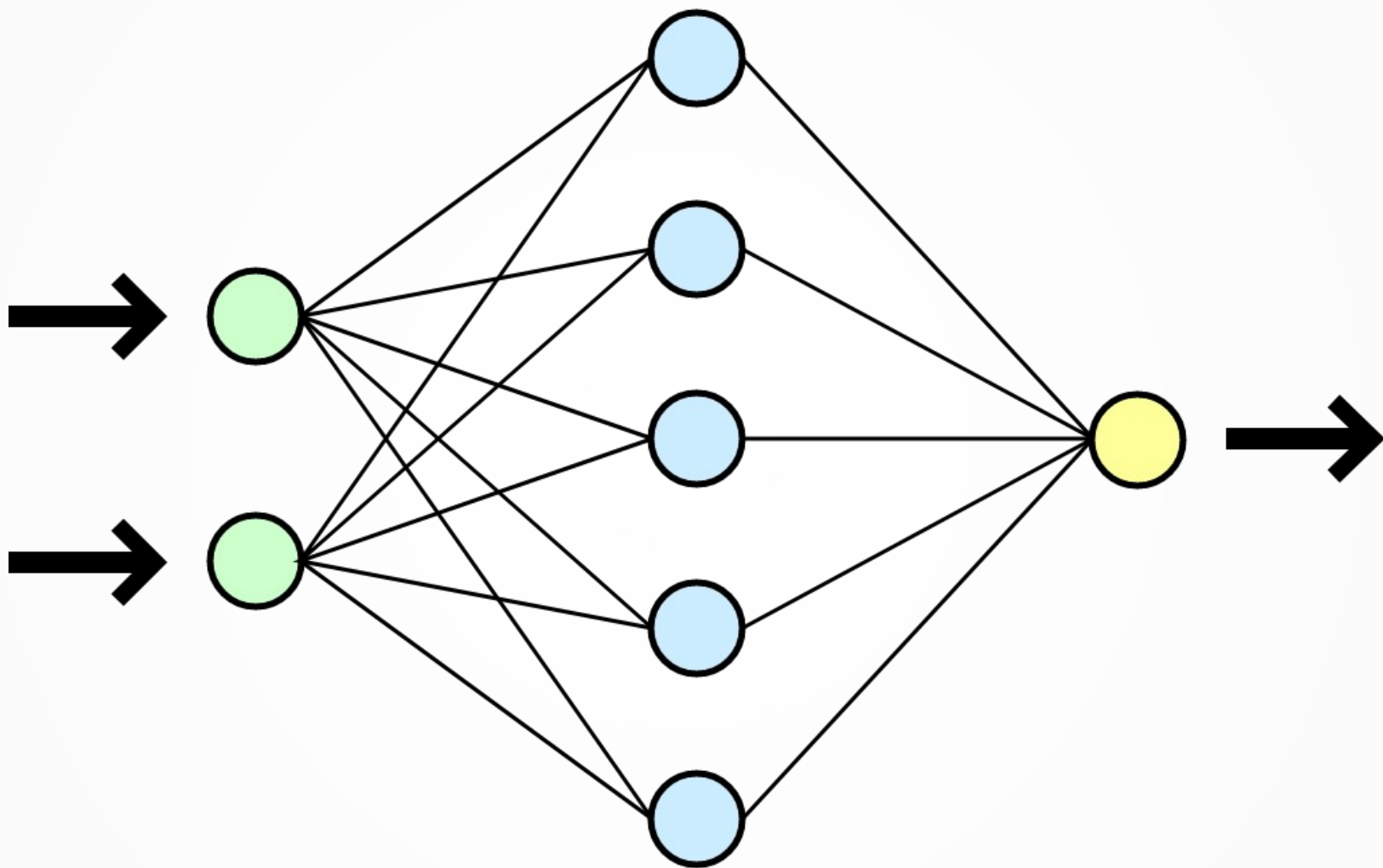
# Literature and Sources



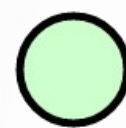
# Sources

- Deep Learning by Andrew W. Trask
  - very practical approach
- Deep Learning by Goodfellow, Bengio, Curville
  - A lot of theoretical background, [online available](#)
- Stanford CS224: Natural Language Processing with Deep Learning
  - [public available](#) course, with videos from all lectures

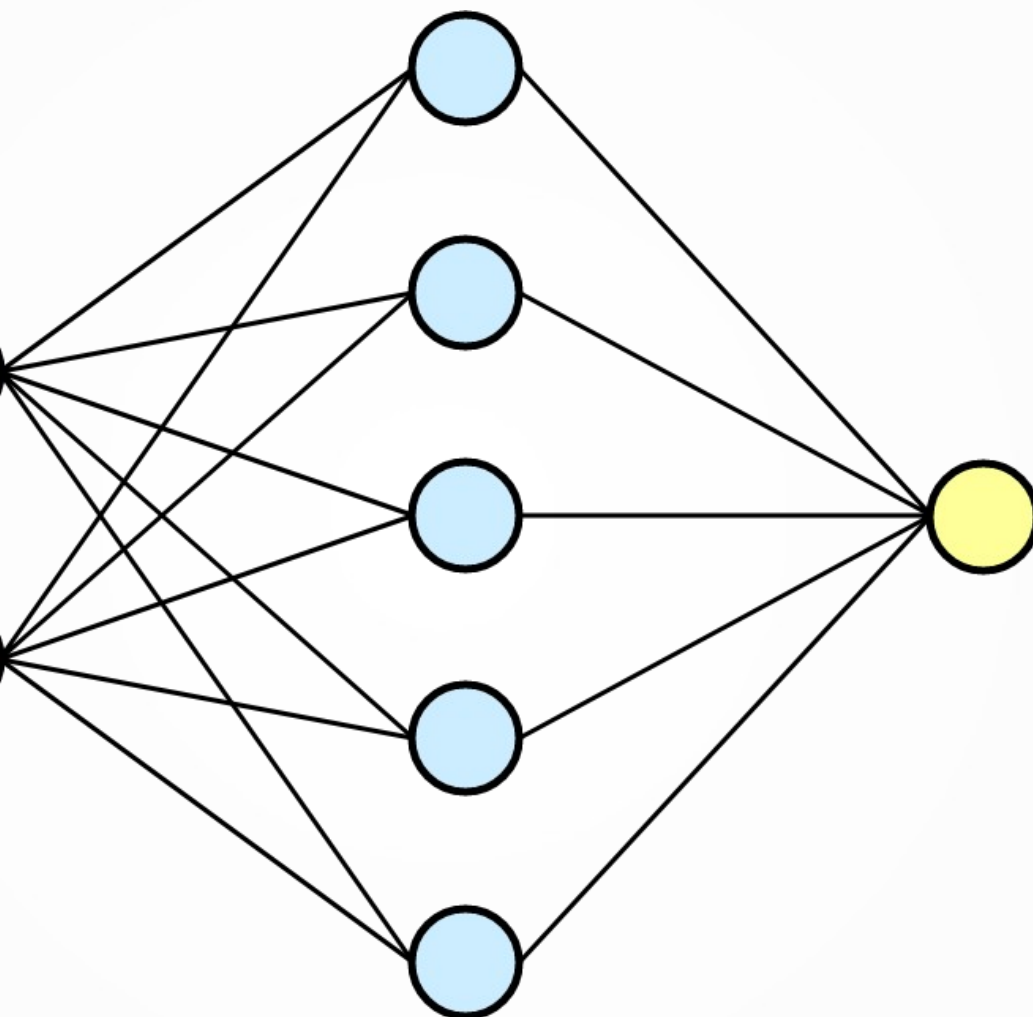
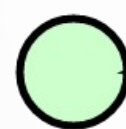
# **Text and Neural Networks**



O

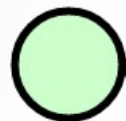


K

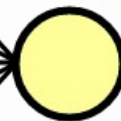
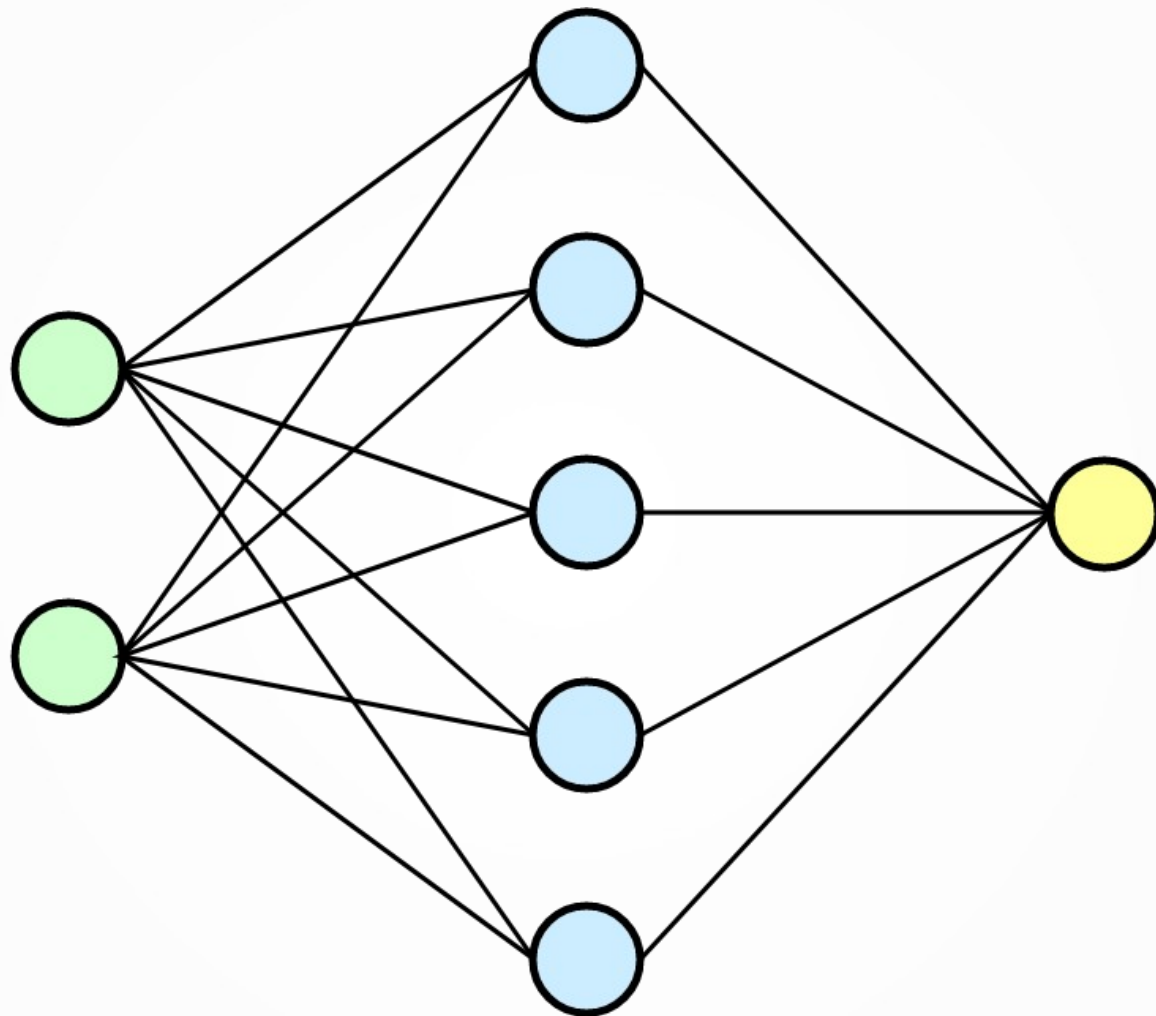
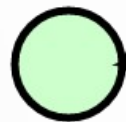


1

n



o



0



# **What about:**

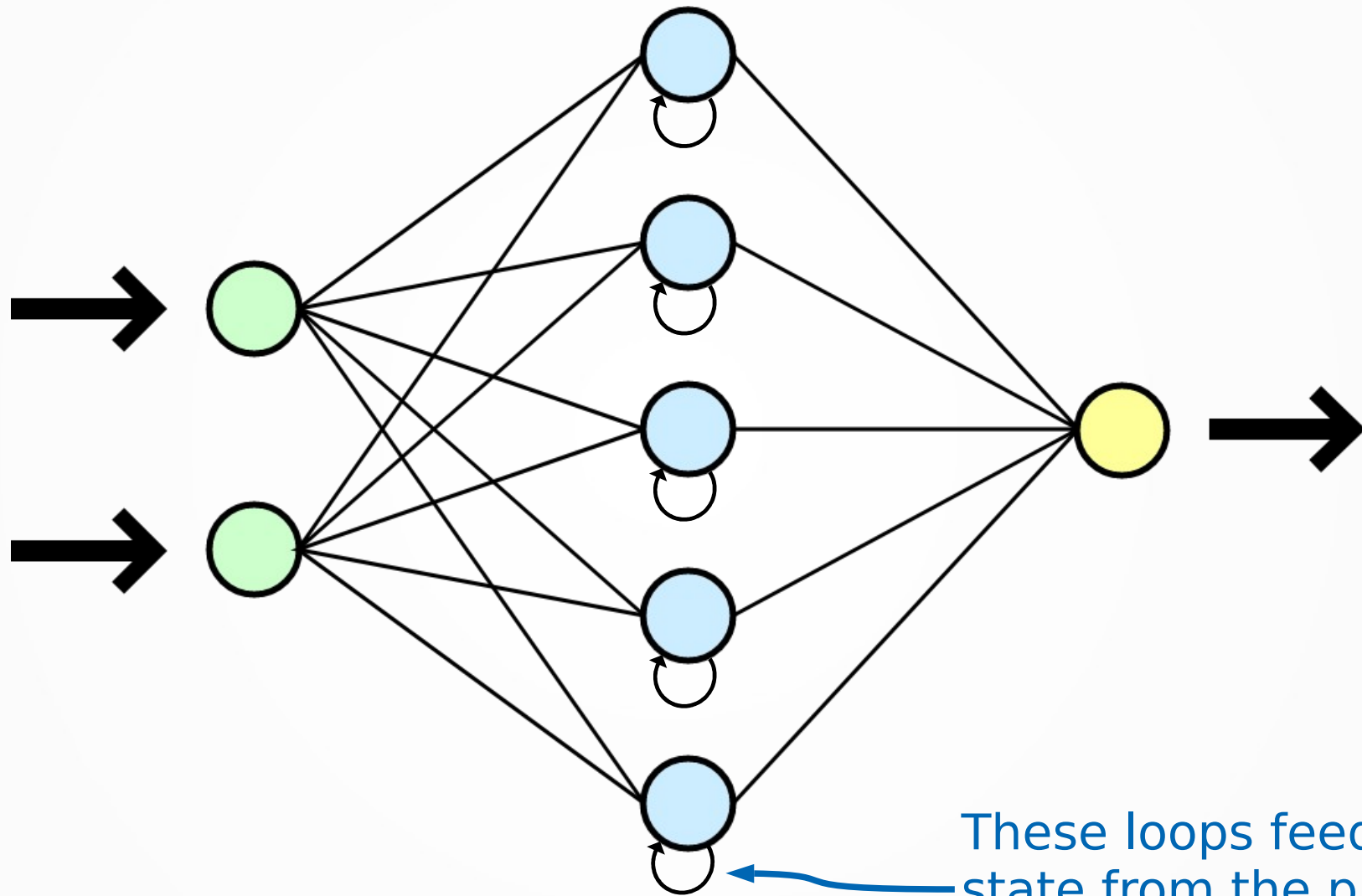
**Not  
Never  
Nope  
yes**

**...**

**We need to handle inputs of  
arbitrary length.**

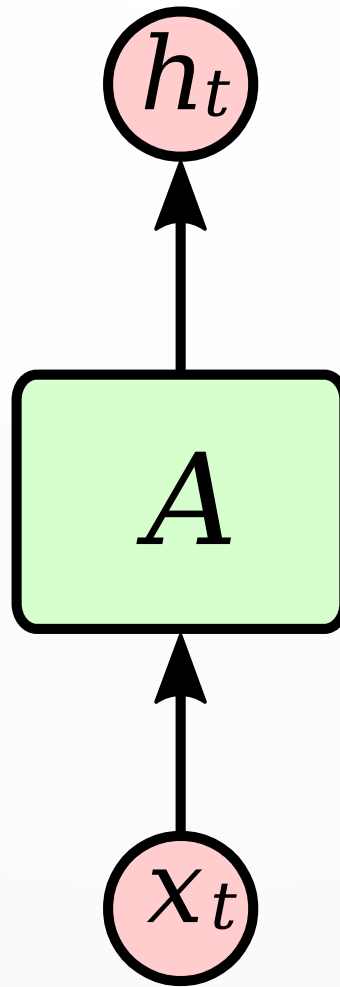
# Recurrent Neural Networks

RNNs are networks with loops, allowing information to persist. [Rummelhart et al. 1986]



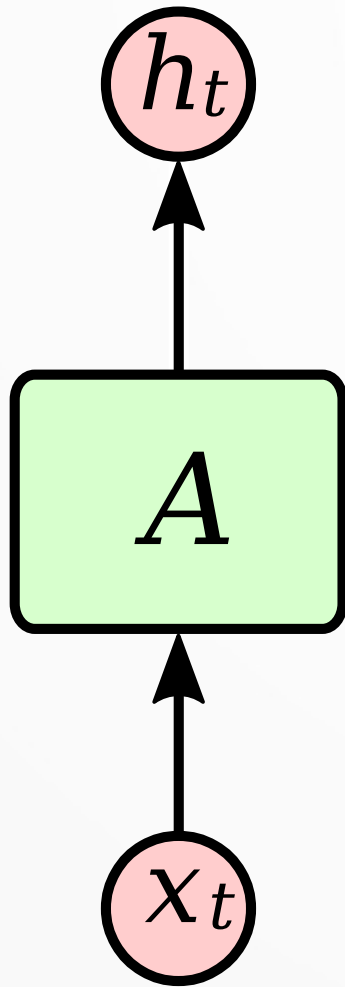
These loops feed in the cell state from the previous timestep.

# Simple Neural Network





# Simple Neural Network



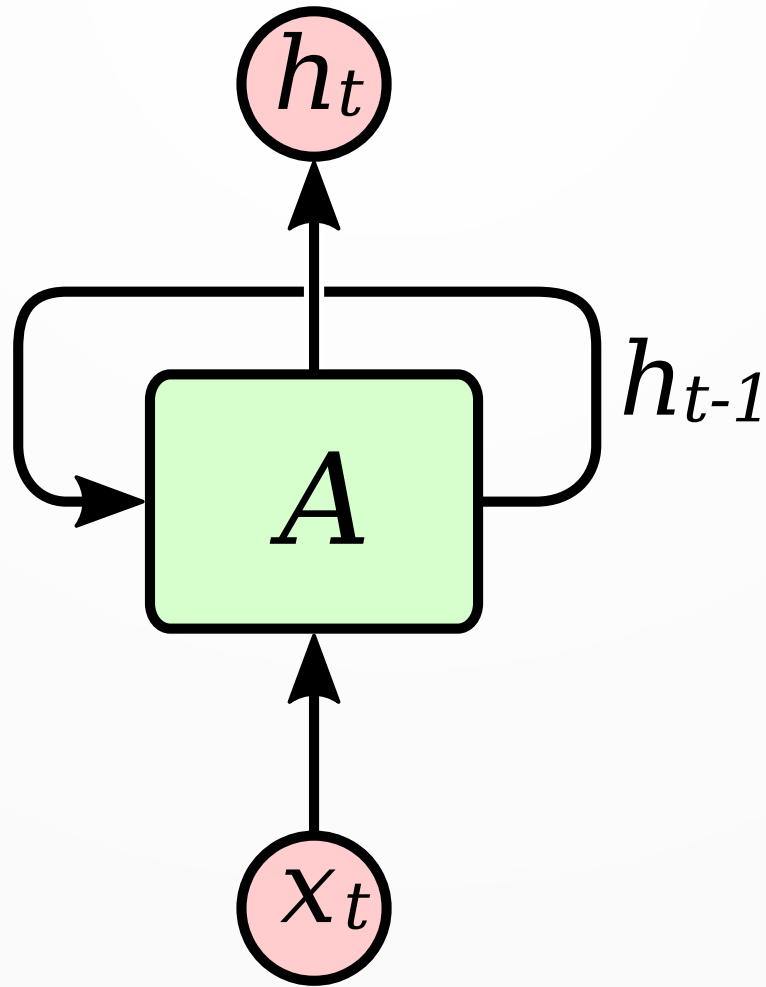
$$\underline{h_t} = \underline{A}(\underline{x_t})$$

new state

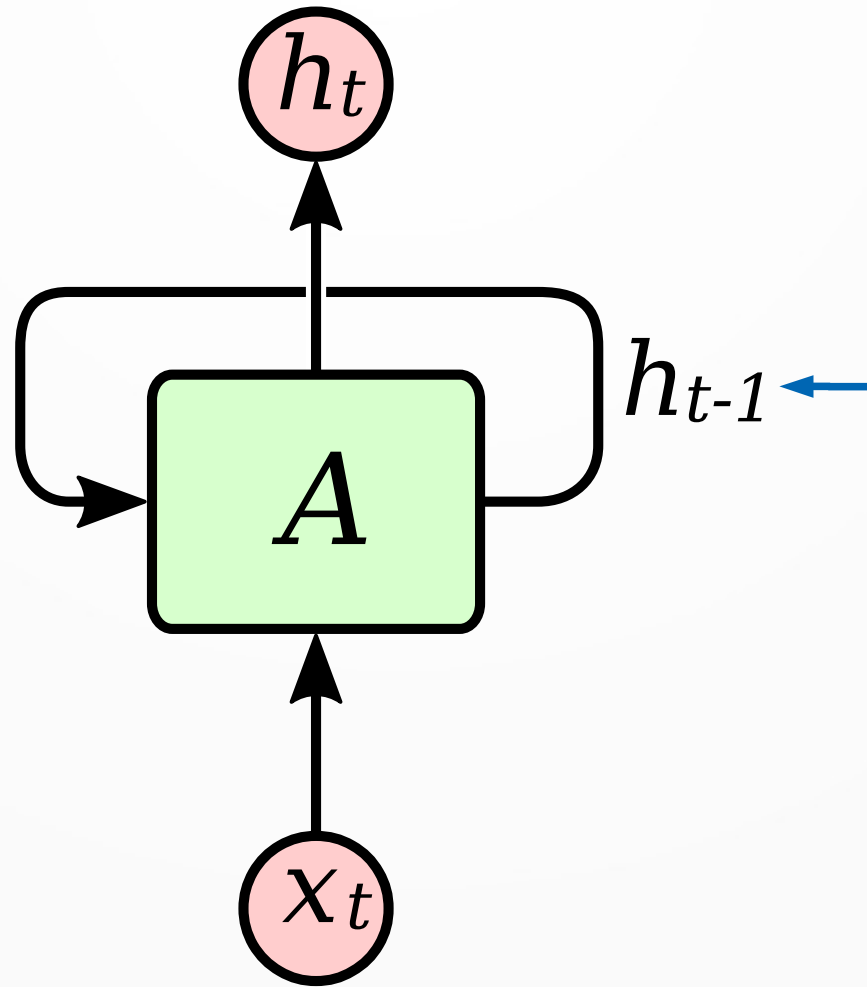
network  
function

input vector

# Recurrent Neural Network

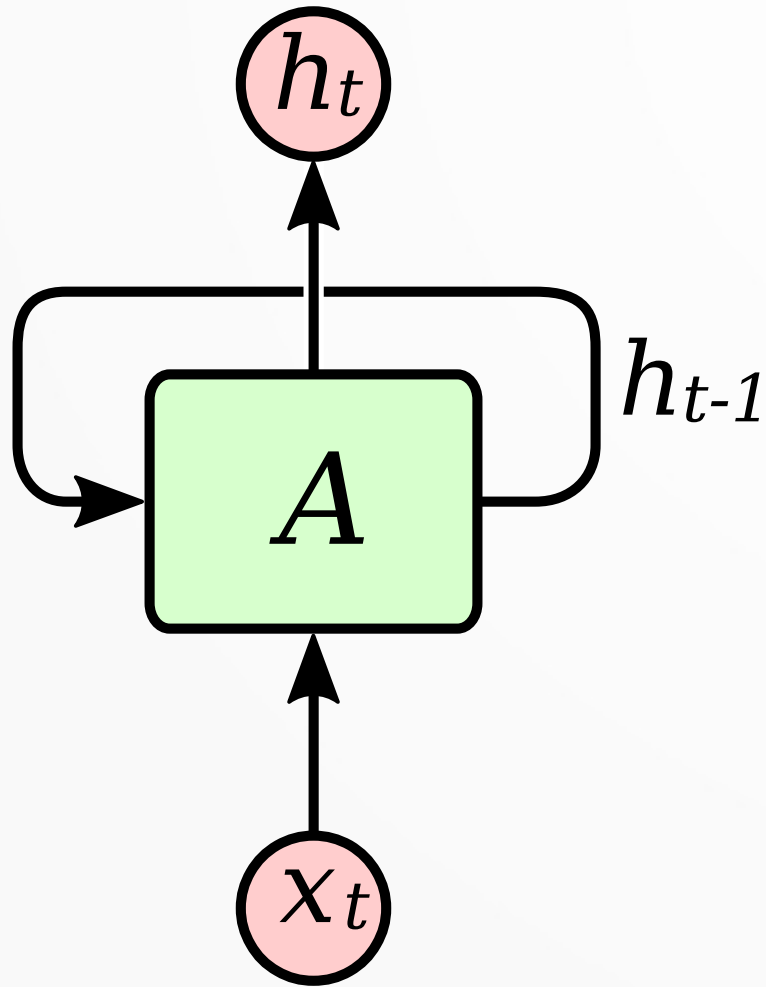


# Recurrent Neural Network



By passing the previous hidden state, the network can keep an „internal state“ as the input sequence is processed.

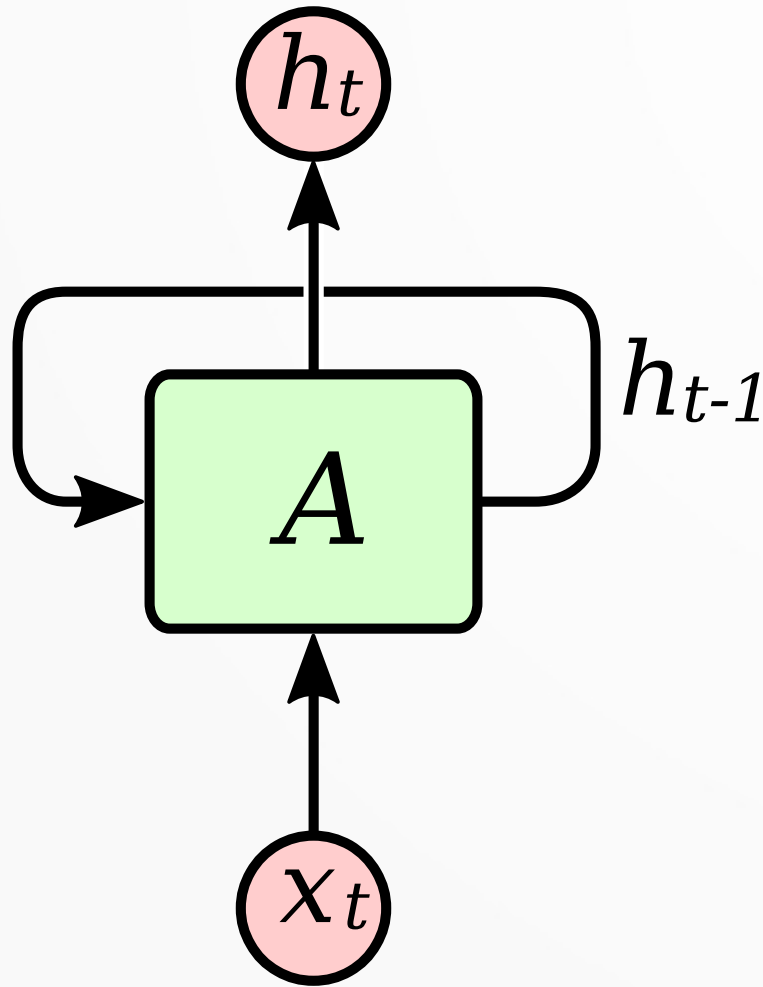
# Recurrent Neural Network



$$\underline{h_t} = \underline{A}(\underline{h_{t-1}}, \underline{x_t})$$

new state      network function      previous state      input vector

# Recurrent Neural Network

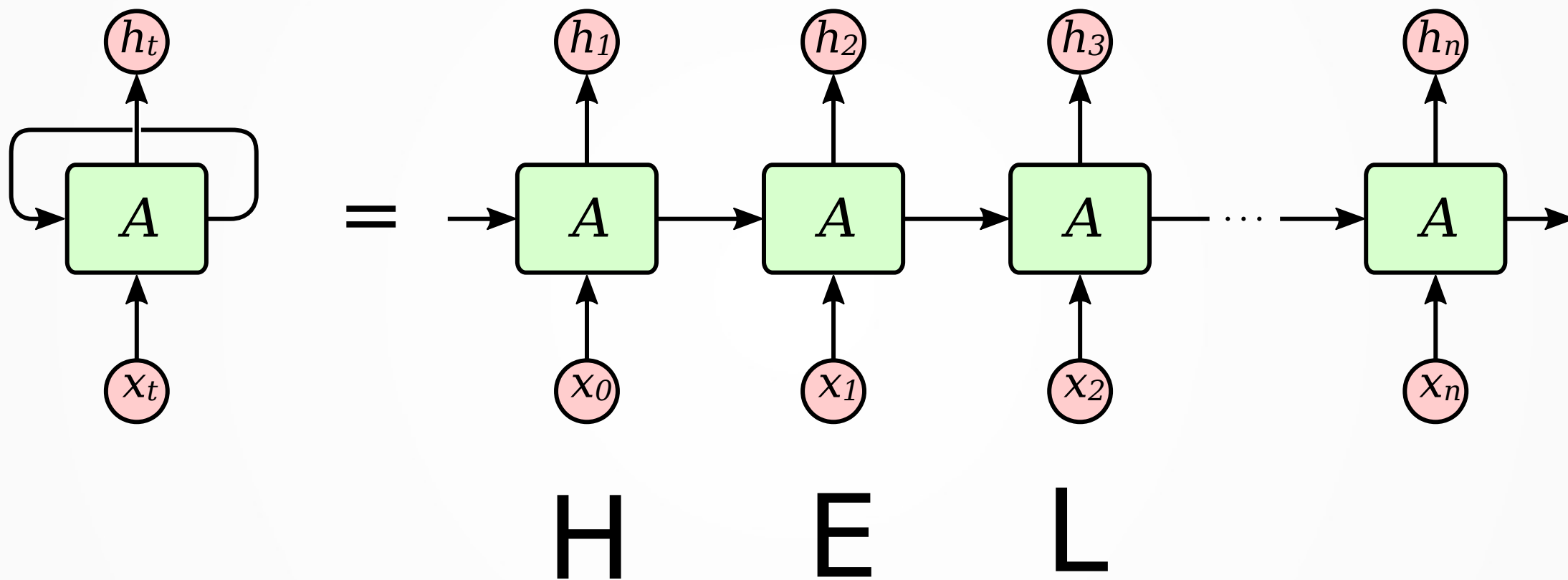


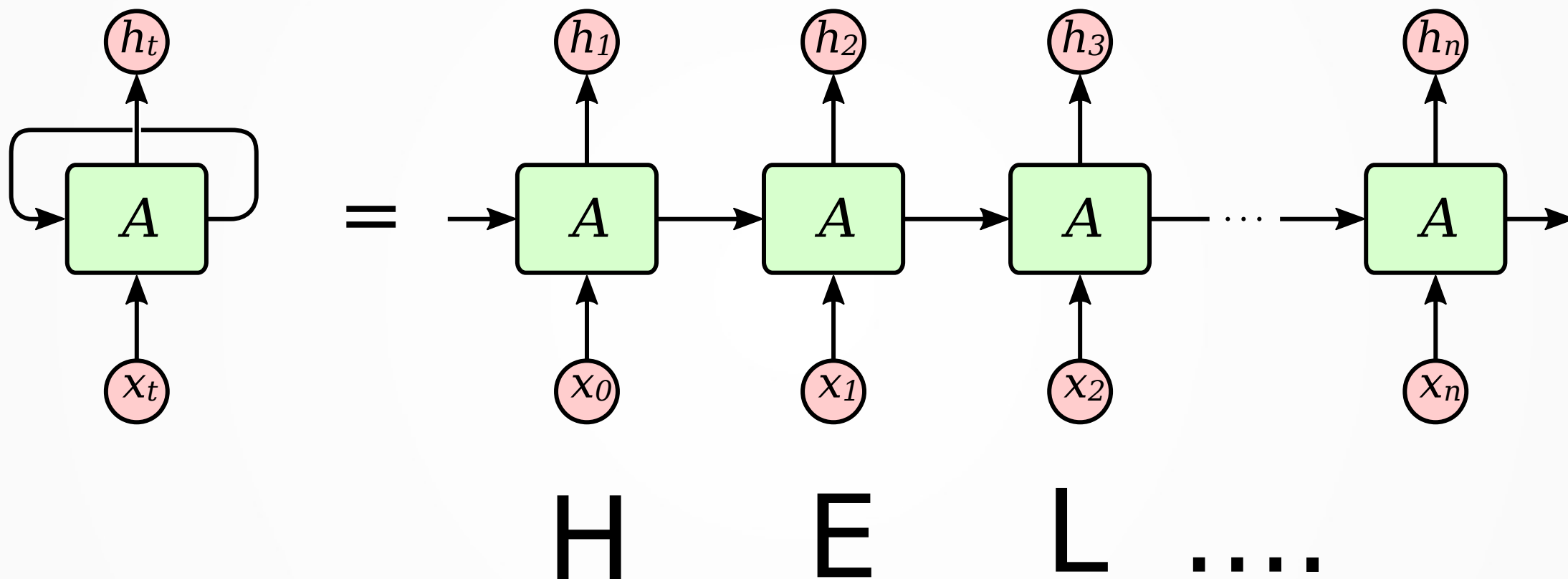
$$\underline{h_t} = \underline{A}(\underline{h_{t-1}}, \underline{x_t})$$

new state      network function      previous state      input vector

We can process a sequence by recursively applying this formula.

# Unfolding in time

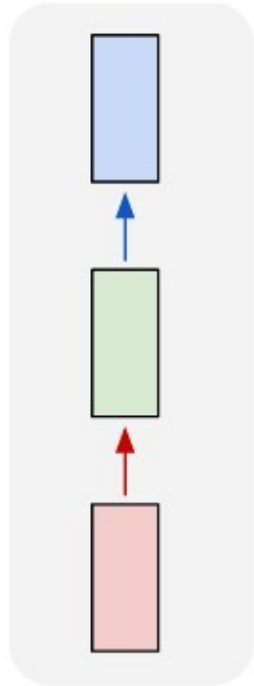






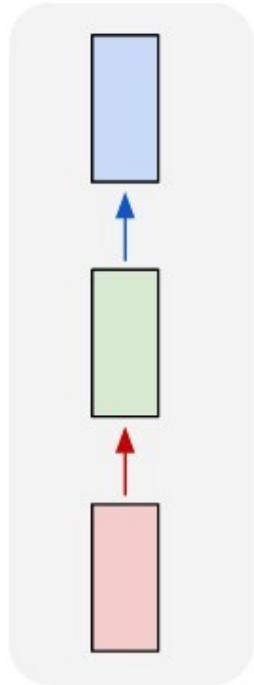
# **Network Architectures**

one to one

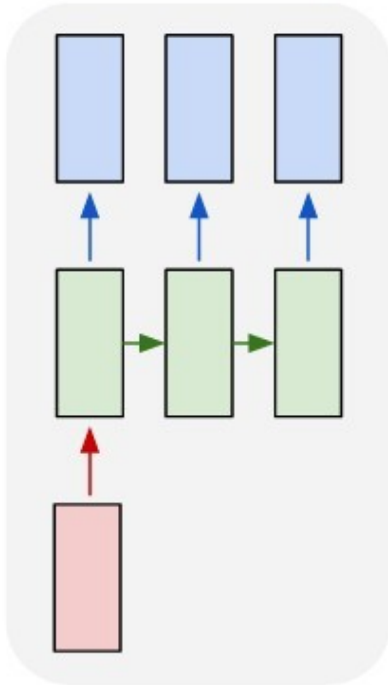


## Simple Neural Network

one to one



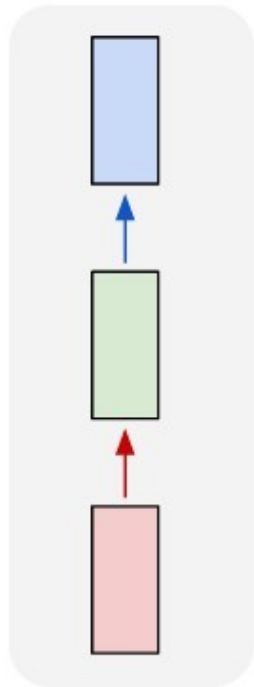
one to many



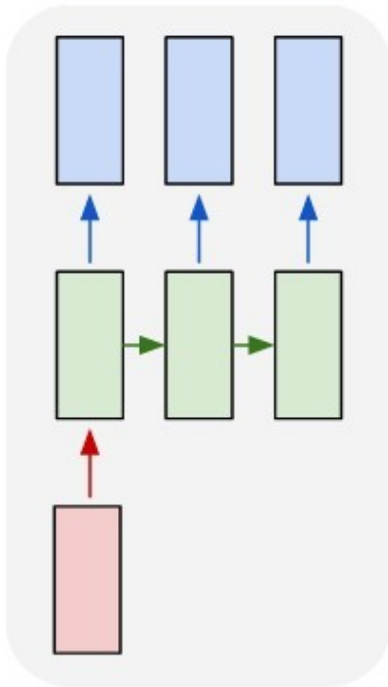
## **image captioning**

Image to a sequence of words

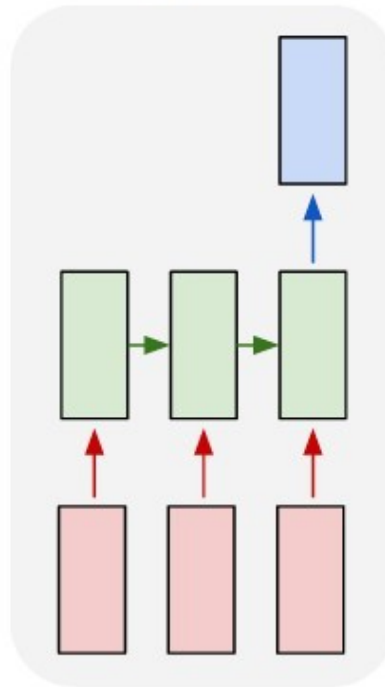
one to one



one to many

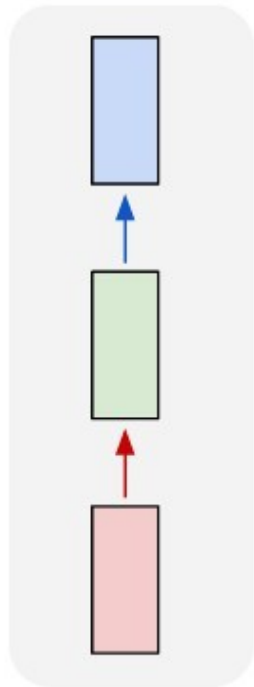


many to one

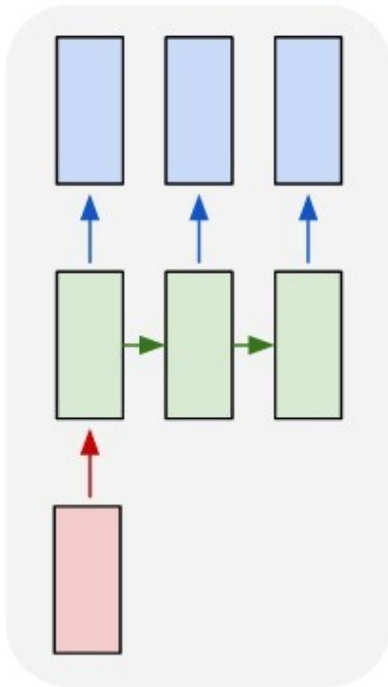


**classification**  
sequence of words to a class

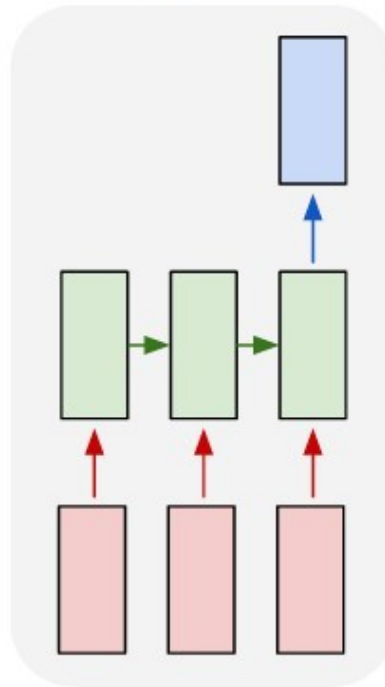
one to one



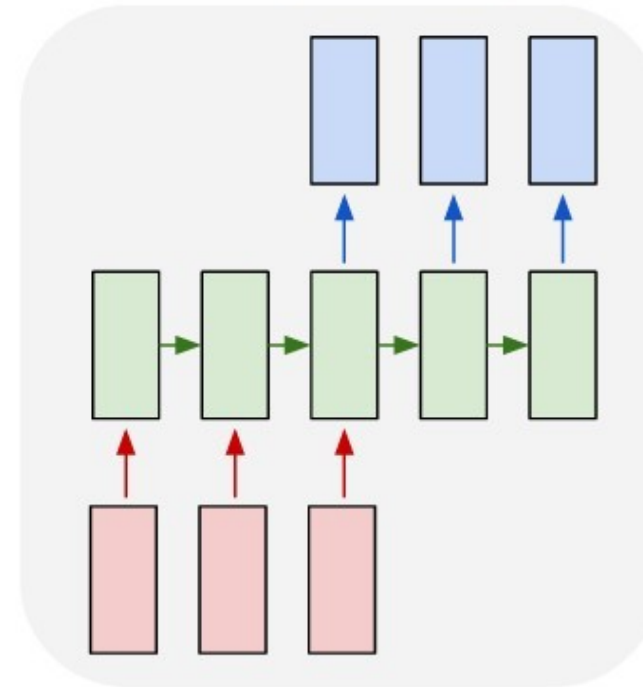
one to many



many to one



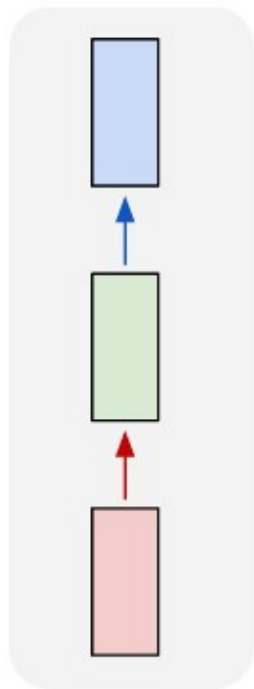
many to many



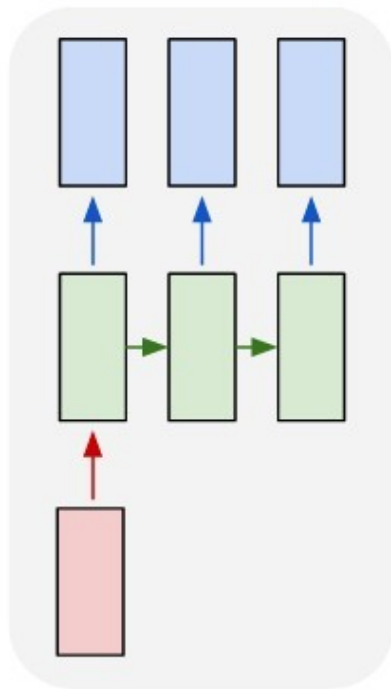
## **machine translation**

sequence of words to sequence of words

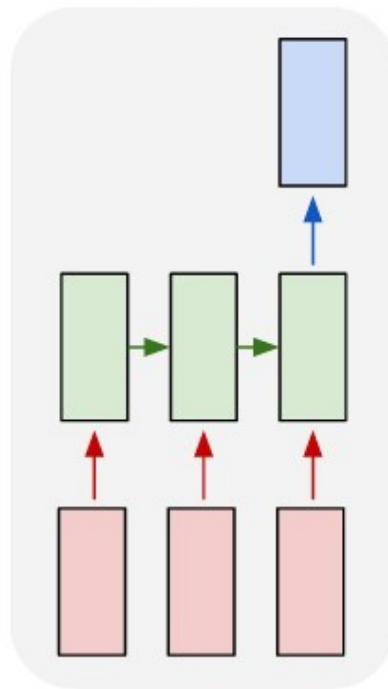
one to one



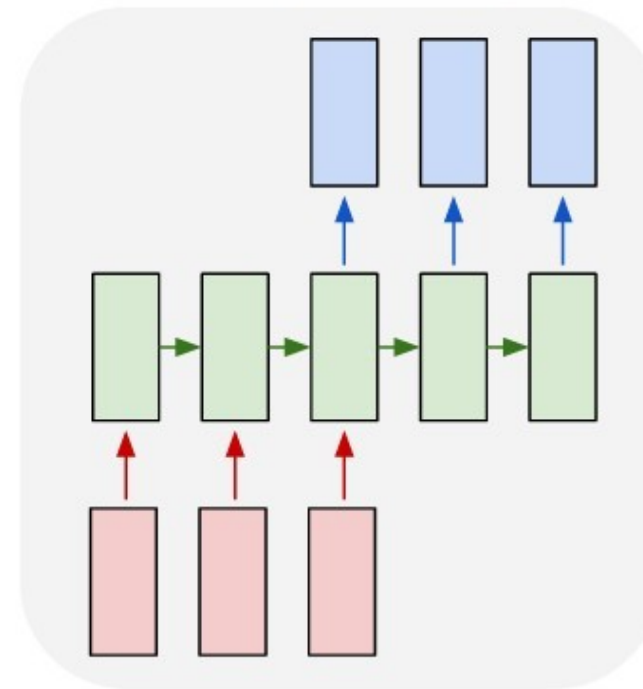
one to many



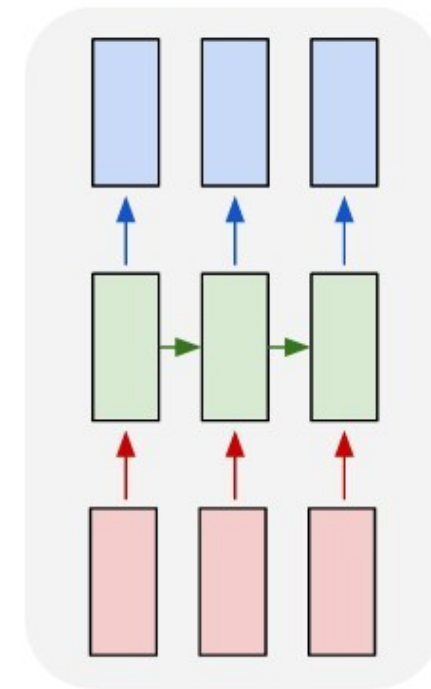
many to one



many to many



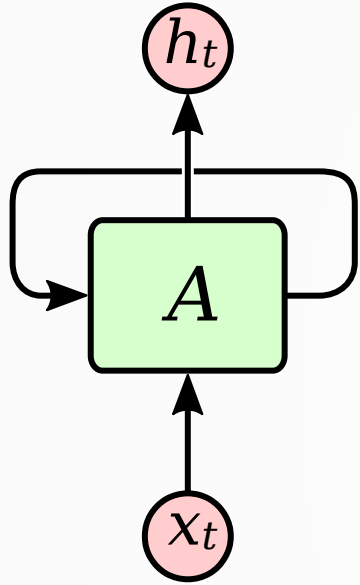
many to many



## Video classification

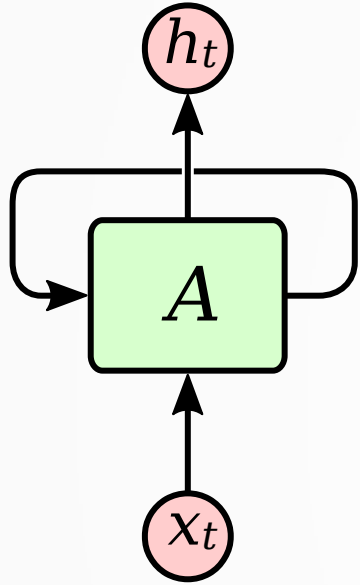
A list of frames to a list of classes

# Backpropagation

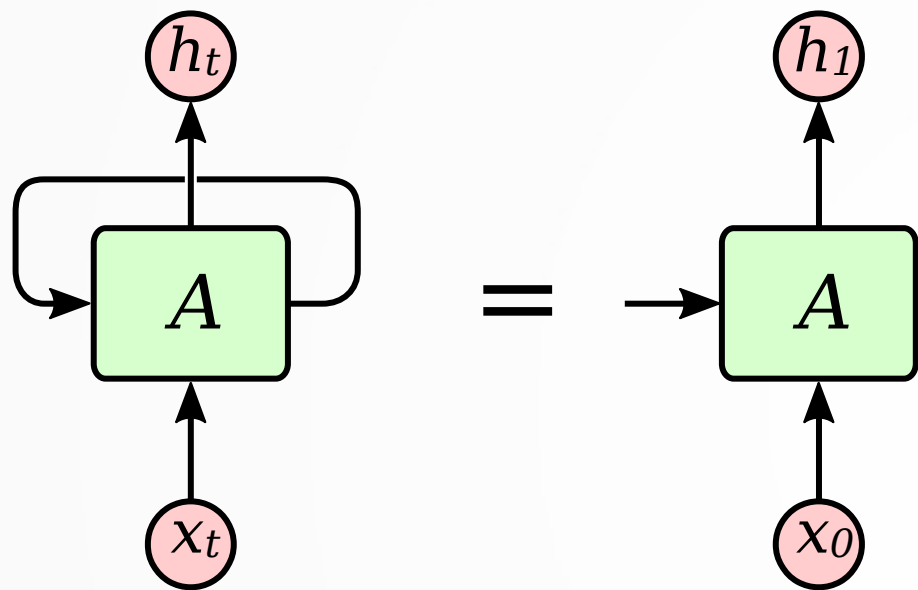


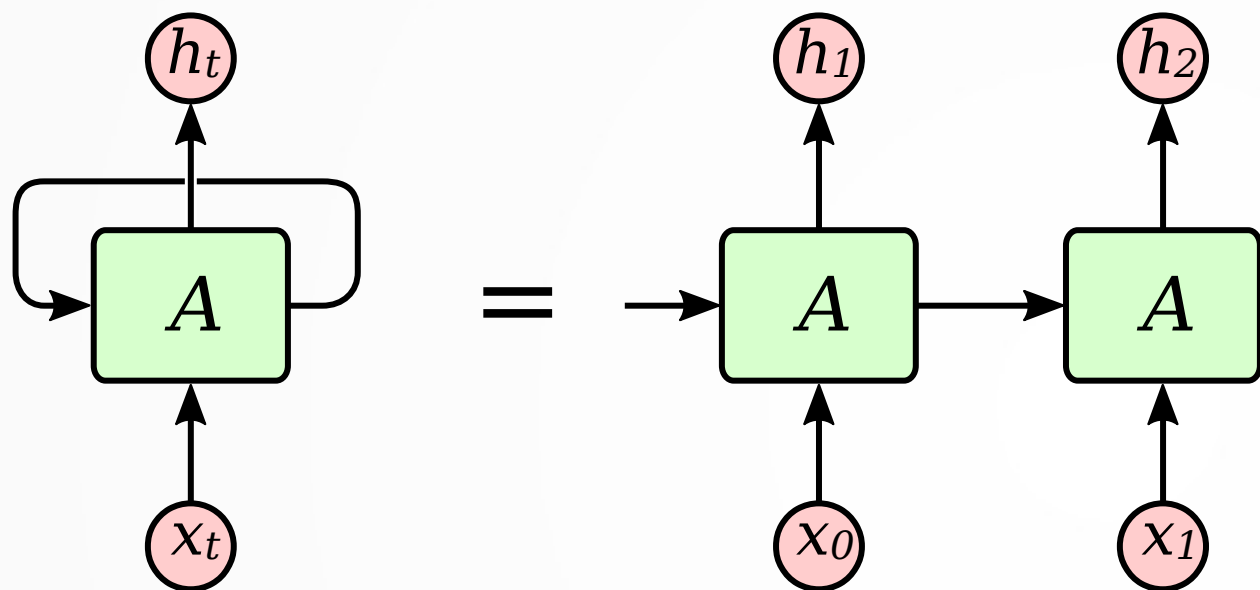
**How can we apply  
backpropagation to nets  
with loops?**



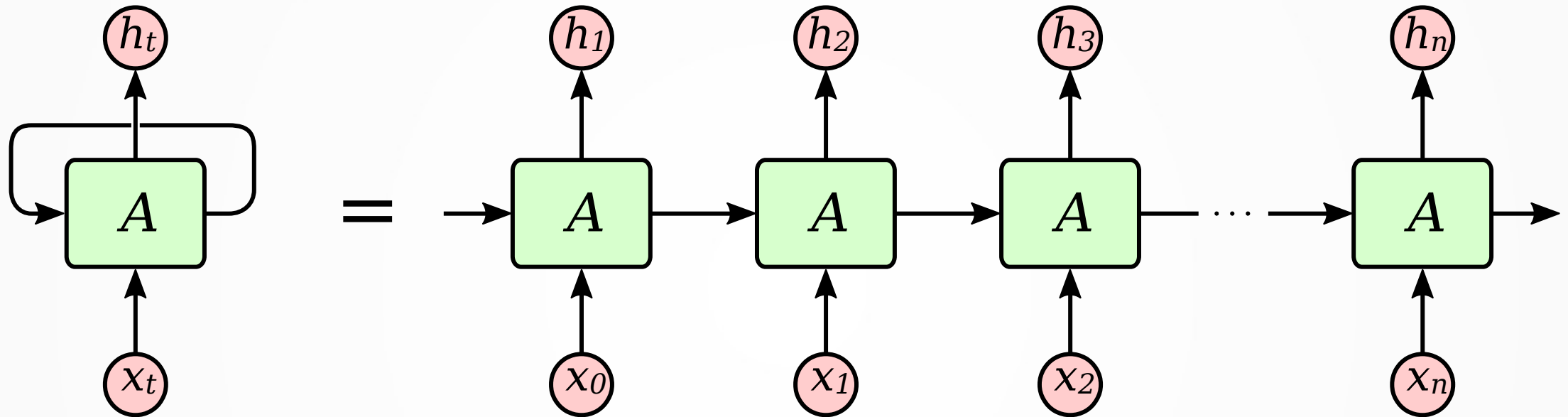


**with backpropagation  
through time (BPTT)**

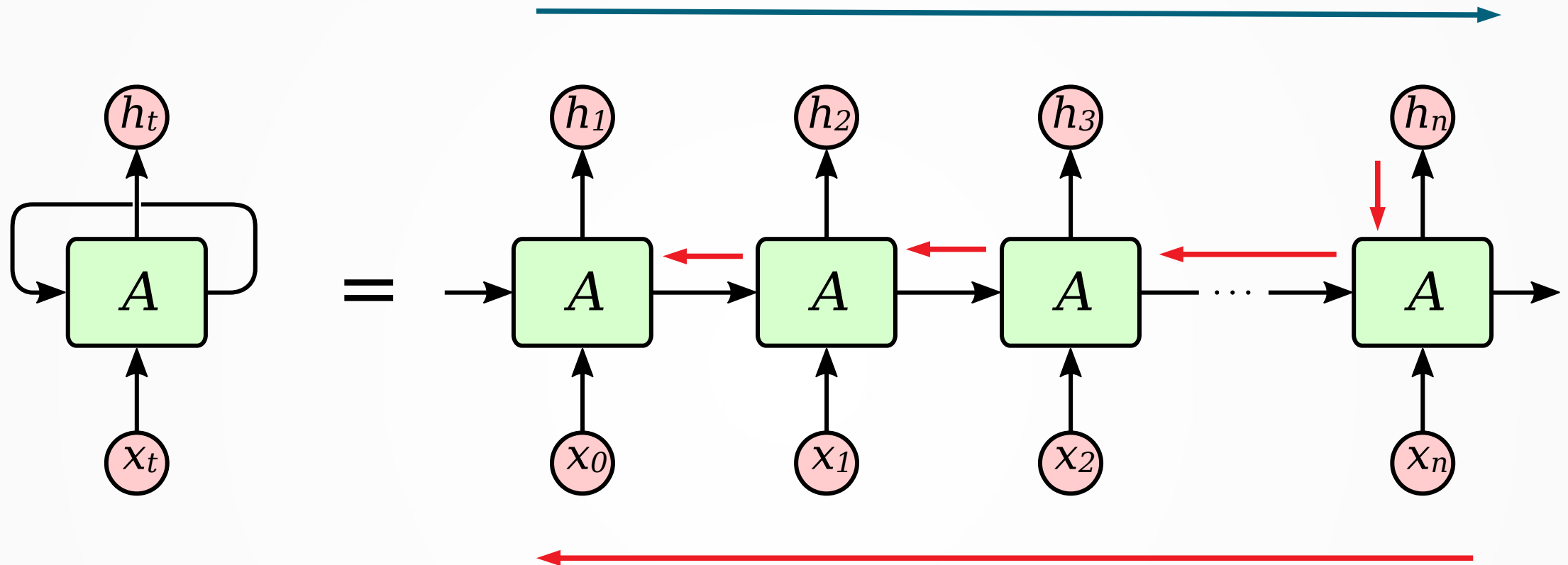




1. forward-propagate the inputs over the unfolded network

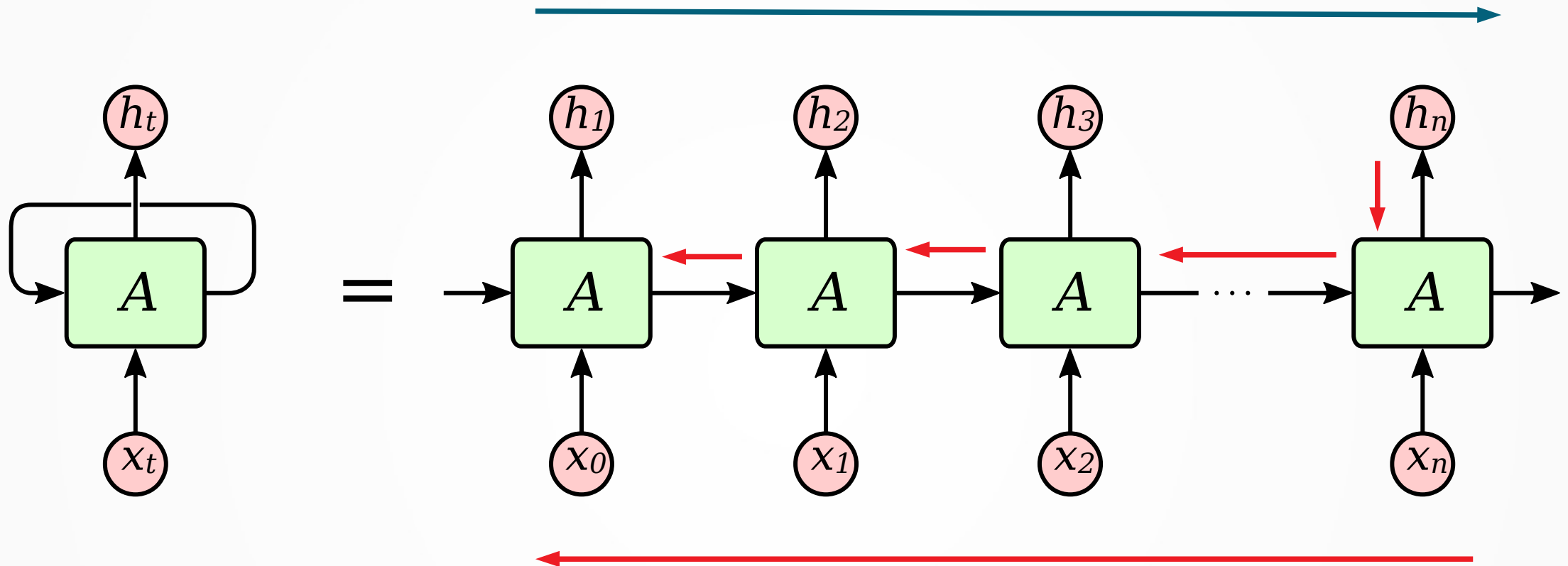


1. forward-propagate the inputs over the unfolded network



2. back-propagate the error, back across the unfolded network

1. forward-propagate the inputs over the unfolded network



2. back-propagate the error, back across the unfolded network

3. sum the weight changes and update all weights

# Further Information

- A well explained implementation of BPTT can be found [here](#)
- Andrew Ng explaining BPTT

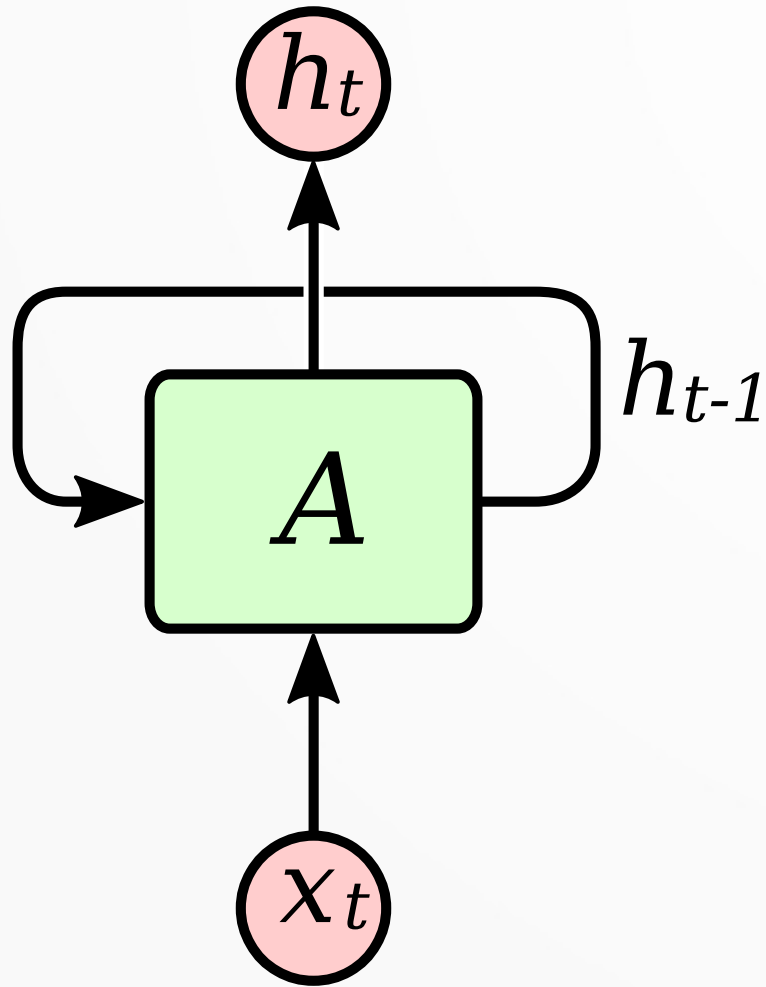
# **The Vanishing Gradient Problem**



# **Learning long-term dependencies with gradient descent is difficult**

Y. Bengio, P. Simard and P. Frasconi in IEEE Transactions on Neural Networks, vol. 5, no. 2, pp. 157-166, March 1994.

# Recurrent Neural Network



$$\underline{h_t} = \underline{A}(\underline{h_{t-1}}, \underline{x_t})$$

new state      network function      previous state      input vector

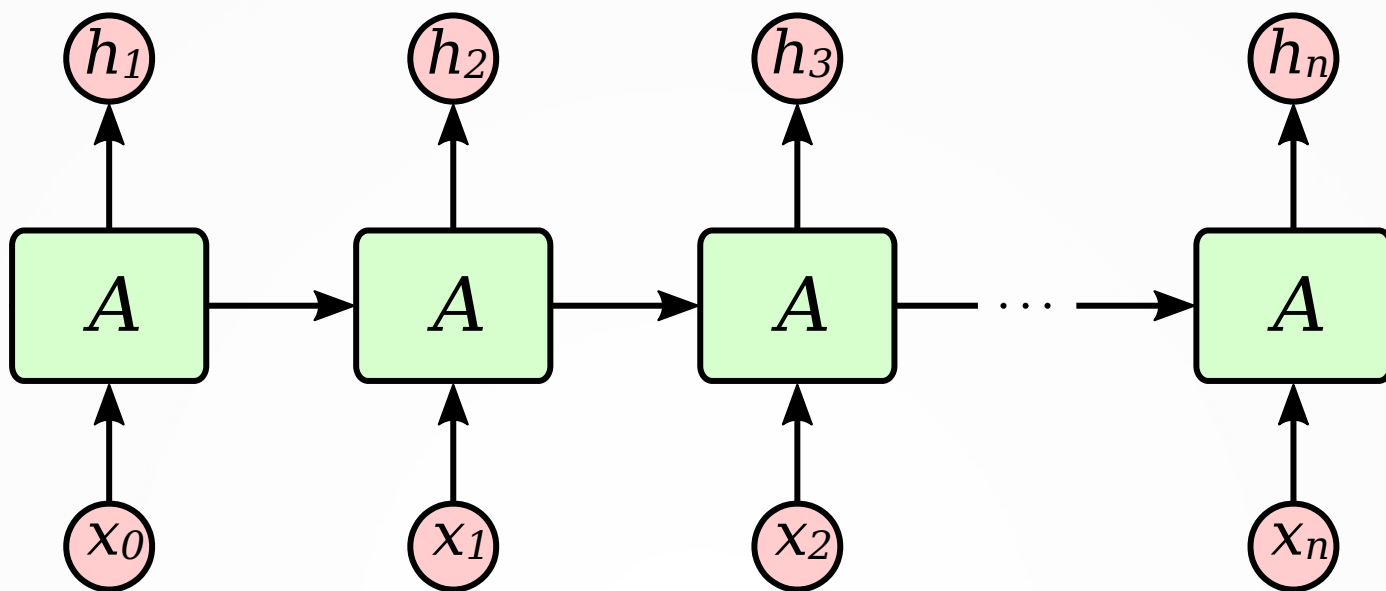
# Recurrent Neural Network

$$\underline{h_t} = \underline{A}(\underline{h_{t-1}}, \underline{x_t})$$

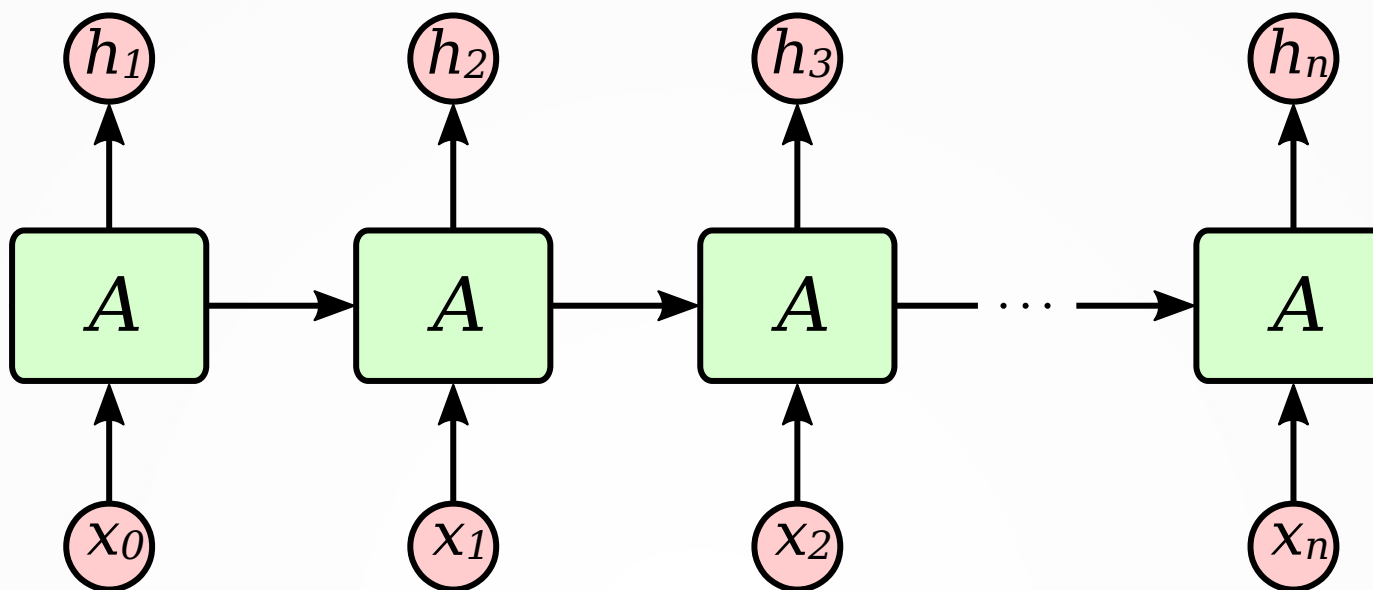
new state      network function      previous state      input vector

A simple example of a network function using tanh

$$\underline{h_t} = \underline{\tanh}(\underline{W_{hh} h_{t-1} + W_{xh} x_t})$$
$$y_t = W_{yh} h_t$$

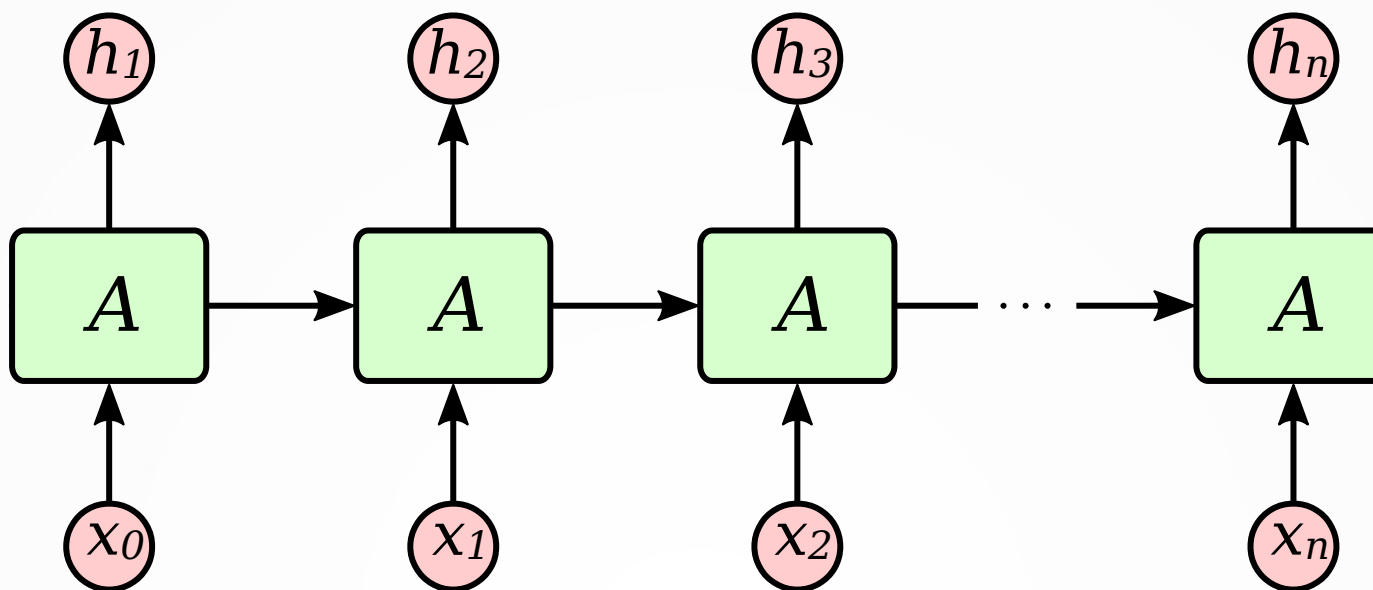


$$h_1 = \tanh(W_{hh}h_0 + W_{xh}x_1)$$



$$h_1 = \tanh(W_{hh} h_0 + W_{xh} x_1)$$

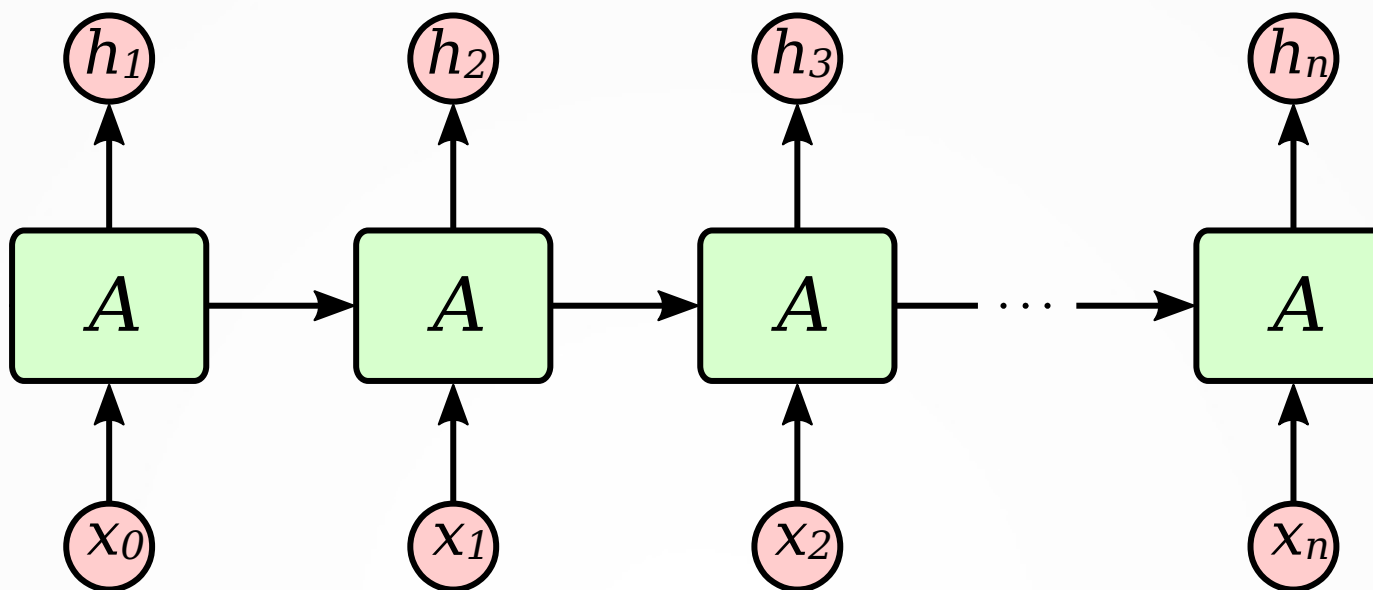
$$h_2 = \tanh(W_{hh} (\tanh(W_{hh} h_0 + W_{xh} x_1)) + W_{xh} x_2)$$



$$h_1 = \tanh(W_{hh} h_0 + W_{xh} x_1)$$

$$h_2 = \tanh(W_{hh}(\tanh(W_{hh} h_0 + W_{xh} x_1)) + W_{xh} x_2)$$

$$h_3 = \tanh(W_{hh}(\tanh(W_{hh}(\tanh(W_{hh} h_0 + W_{xh} x_1)) + W_{xh} x_2)) + W_{xh} x_3)$$



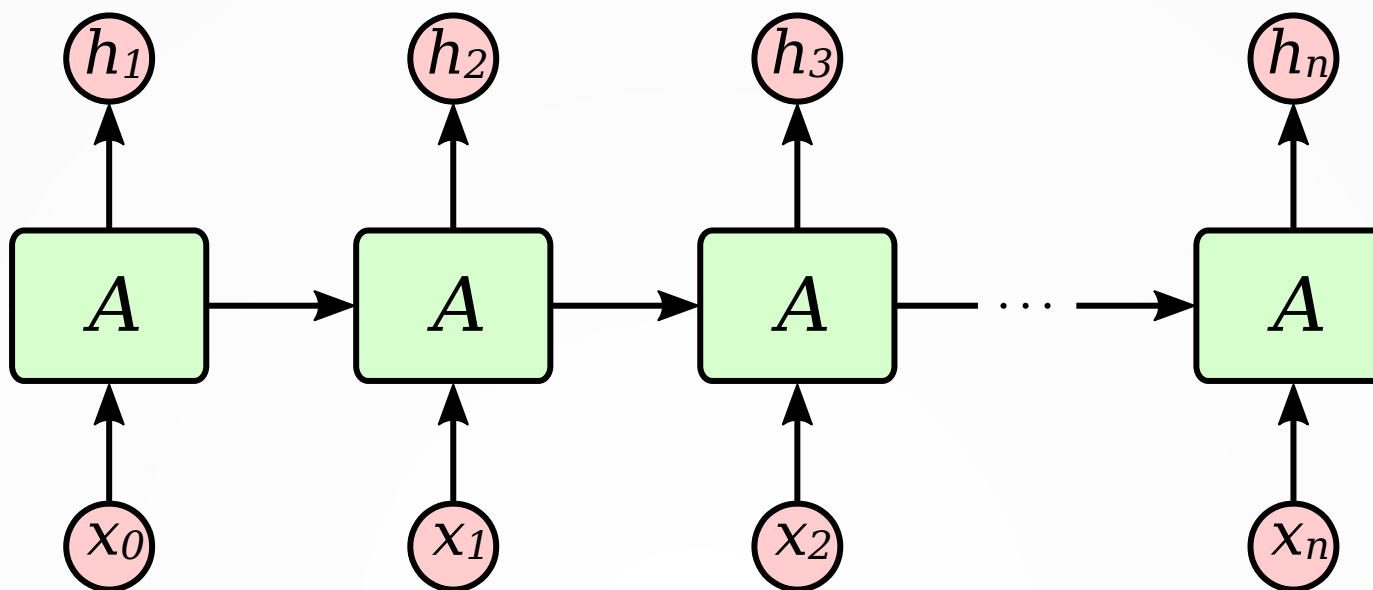
$$h_1 = \tanh(W_{hh} h_0 + W_{xh} x_1)$$

$$h_2 = \tanh(W_{hh}(\tanh(W_{hh} h_0 + W_{xh} x_1)) + W_{xh} x_2)$$

$$h_3 = \tanh(W_{hh}(\tanh(W_{hh}(\tanh(W_{hh} h_0 + W_{xh} x_1)) + W_{xh} x_2)) + W_{xh} x_3)$$

$$h_4 = \tanh(W_{hh}(\tanh(W_{hh}(\tanh(W_{hh}(\tanh(W_{hh} h_0 + W_{xh} x_1)) + W_{xh} x_2)) + W_{xh} x_3)) + W_{xh} x_4)$$





$$h_1 = \tanh(W_{hh} h_0 + W_{xh} x_1)$$

$$h_2 = \tanh(W_{hh}(\tanh(W_{hh} h_0 + W_{xh} x_1)) + W_{xh} x_2)$$

$$h_3 = \tanh(W_{hh}(\tanh(W_{hh}(\tanh(W_{hh} h_0 + W_{xh} x_1)) + W_{xh} x_2)) + W_{xh} x_3)$$

$$h_4 = \tanh(W_{hh}(\tanh(W_{hh}(\tanh(W_{hh}(\tanh(W_{hh} h_0 + W_{xh} x_1)) + W_{xh} x_2)) + W_{xh} x_3)) + W_{xh} x_4)$$



Backpropagating this recursive function leads to exploding or vanishing gradients.



# Let's blow up some gradients

Open the [vanishing gradients notebook](#)



# Papers

On the difficulty of training recurrent neural networks

Pascanu, Mikolov and Bengio, 2013

<http://proceedings.mlr.press/v28/pascanu13.pdf>

Learning long-term dependencies with gradient descent is difficult

Bengio, Simard and Frasconi, 1994

<https://ieeexplore.ieee.org/document/279181>

Untersuchungen zu dynamischen neuronalen Netzen

Hochreiter, 1991

<http://people.idsia.ch/~juergen/SeppHochreiter1991ThesisAdvisorSchmidhuber.pdf>

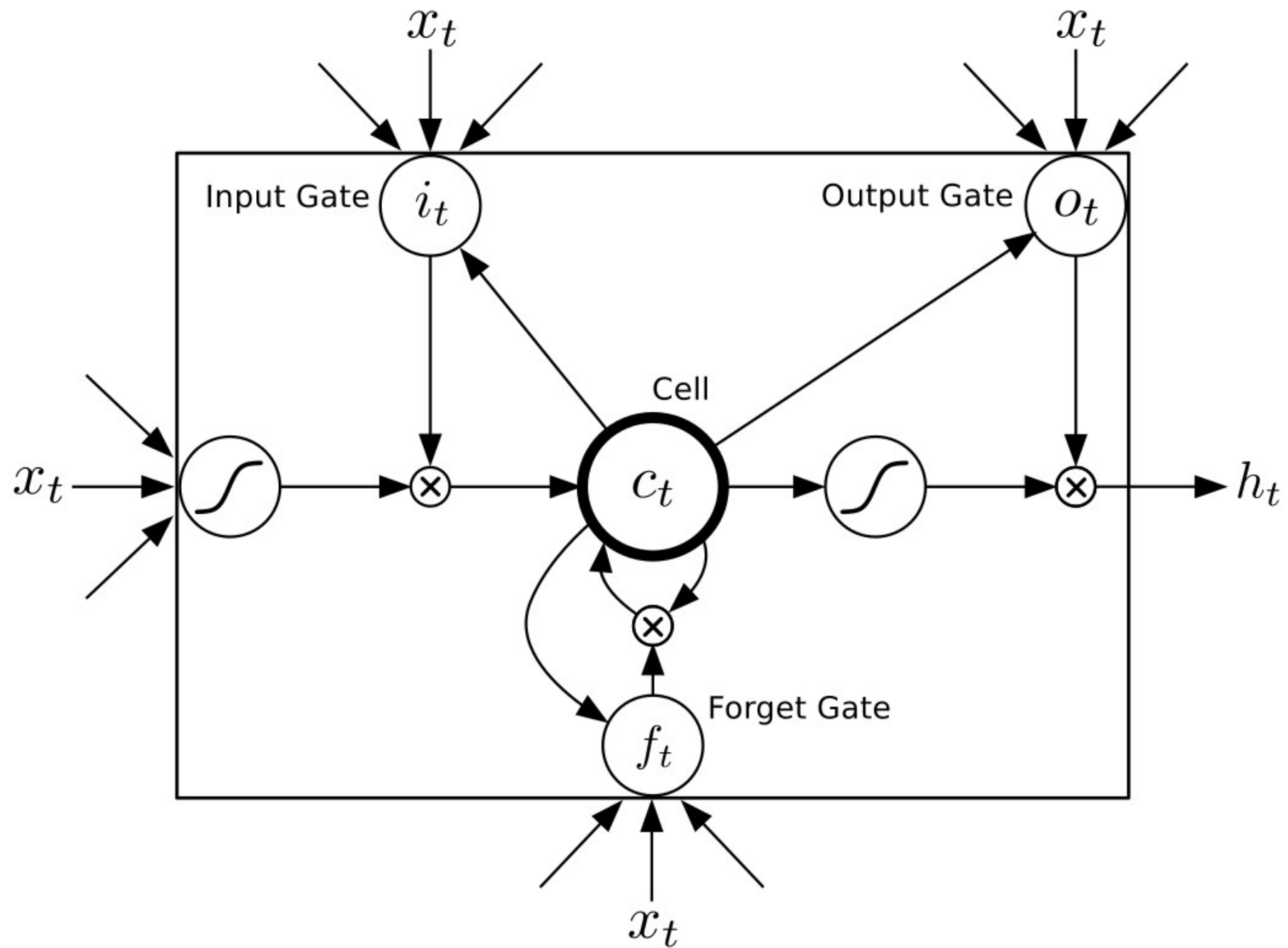
# Solutions for this problem:

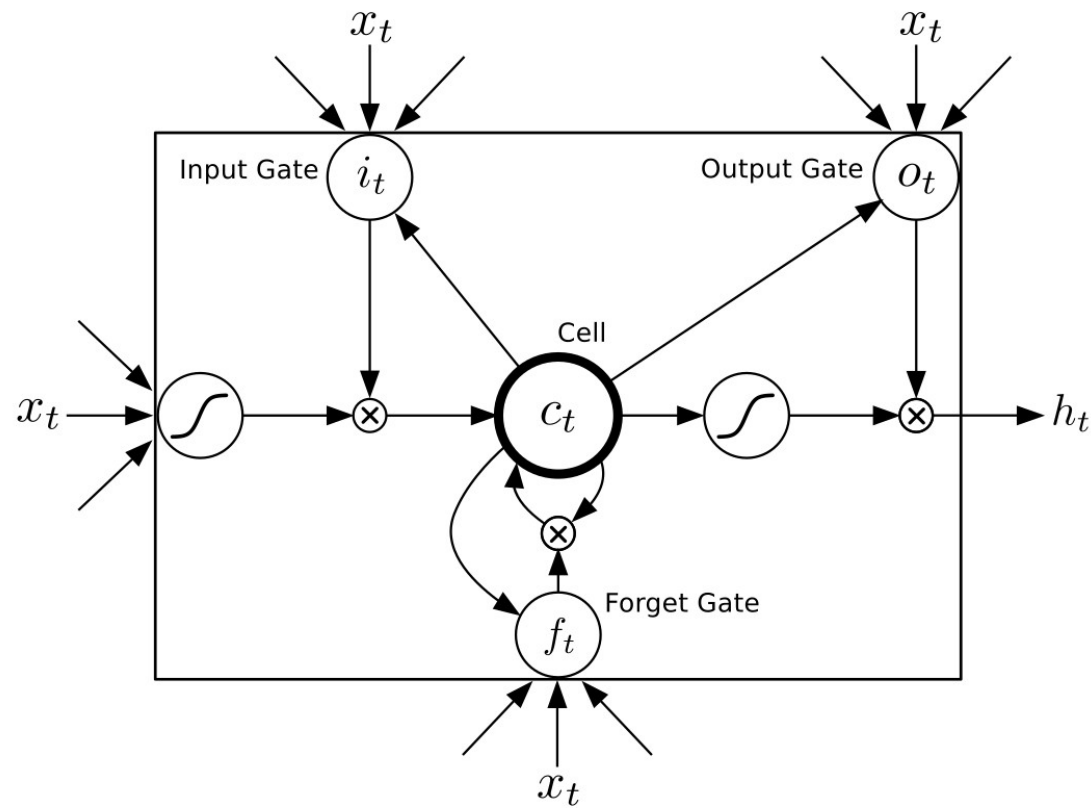
- Limiting the number of past timesteps (Hochreiter, 1991)
- Exploding gradient can be fixed with gradient clipping
- Vanishing gradients can be controlled different architectures (LSTM)
- New: Not using recursions :)

# **LSTM**

**Long short-term memory**

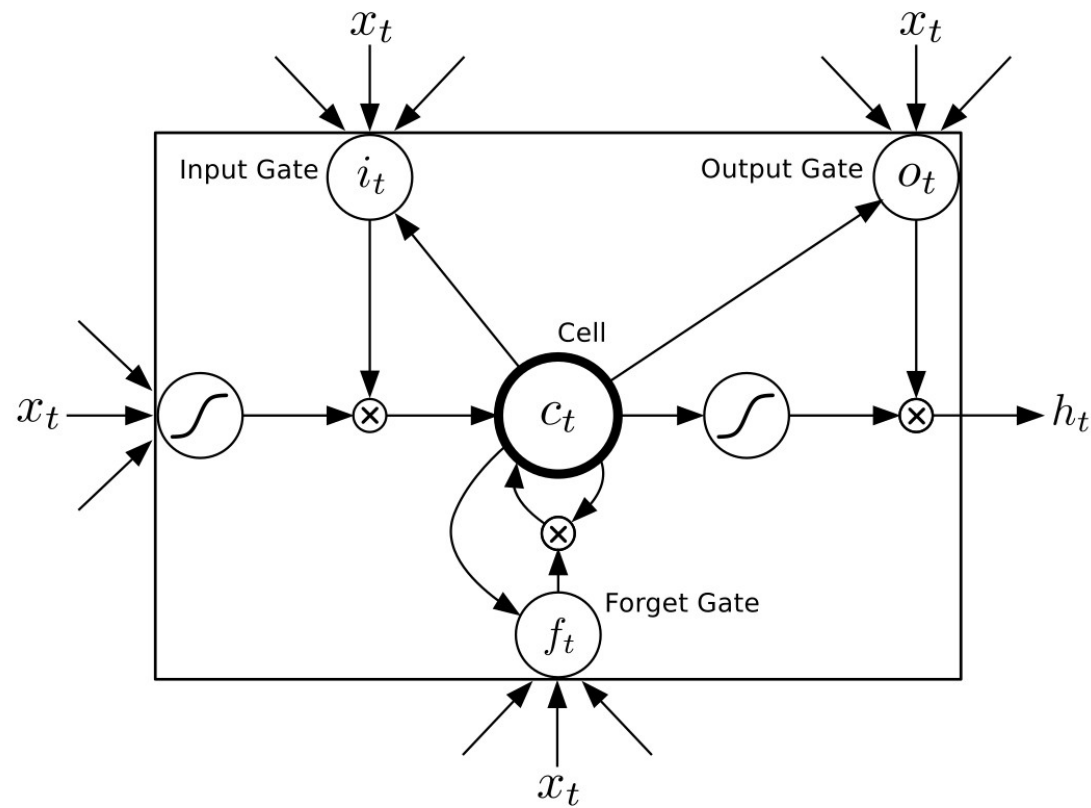
**[Hochreiter et al., 1997]**





$$\begin{aligned}
 i_t &= \sigma(W_{xi}x_t + W_{hi}h_{t-1} + W_{ci}c_{t-1} + b_i) \\
 f_t &= \sigma(W_{xf}x_t + W_{hf}h_{t-1} + W_{cf}c_{t-1} + b_f) \\
 c_t &= f_t c_{t-1} + i_t \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c) \\
 o_t &= \sigma(W_{xo}x_t + W_{ho}h_{t-1} + W_{co}c_t + b_o) \\
 h_t &= o_t \tanh(c_t)
 \end{aligned}$$

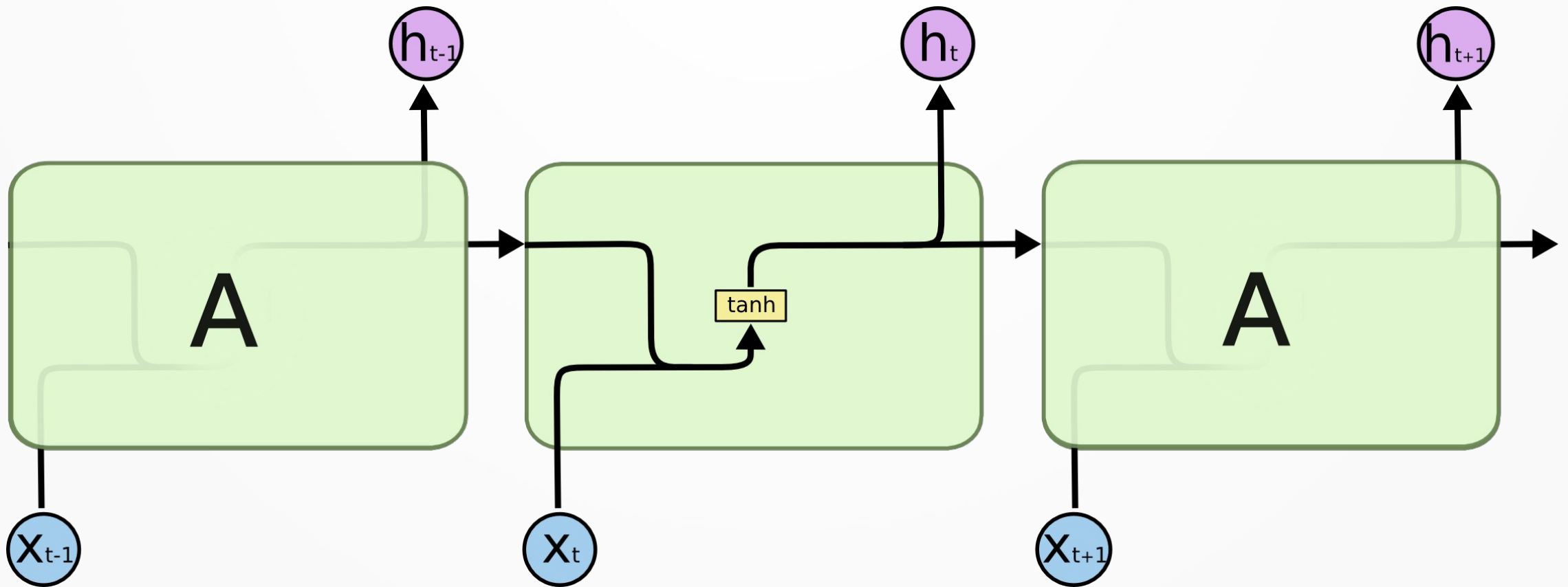




The sigmoid function outputs a number between 0 and 1

$$\begin{aligned}
 i_t &= \sigma(W_{xi}x_t + W_{hi}h_{t-1} + W_{ci}c_{t-1} + b_i) \\
 f_t &= \sigma(W_{xf}x_t + W_{hf}h_{t-1} + W_{cf}c_{t-1} + b_f) \\
 c_t &= f_t c_{t-1} + i_t \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c) \\
 o_t &= \sigma(W_{xo}x_t + W_{ho}h_{t-1} + W_{co}c_t + b_o) \\
 h_t &= o_t \tanh(c_t)
 \end{aligned}$$

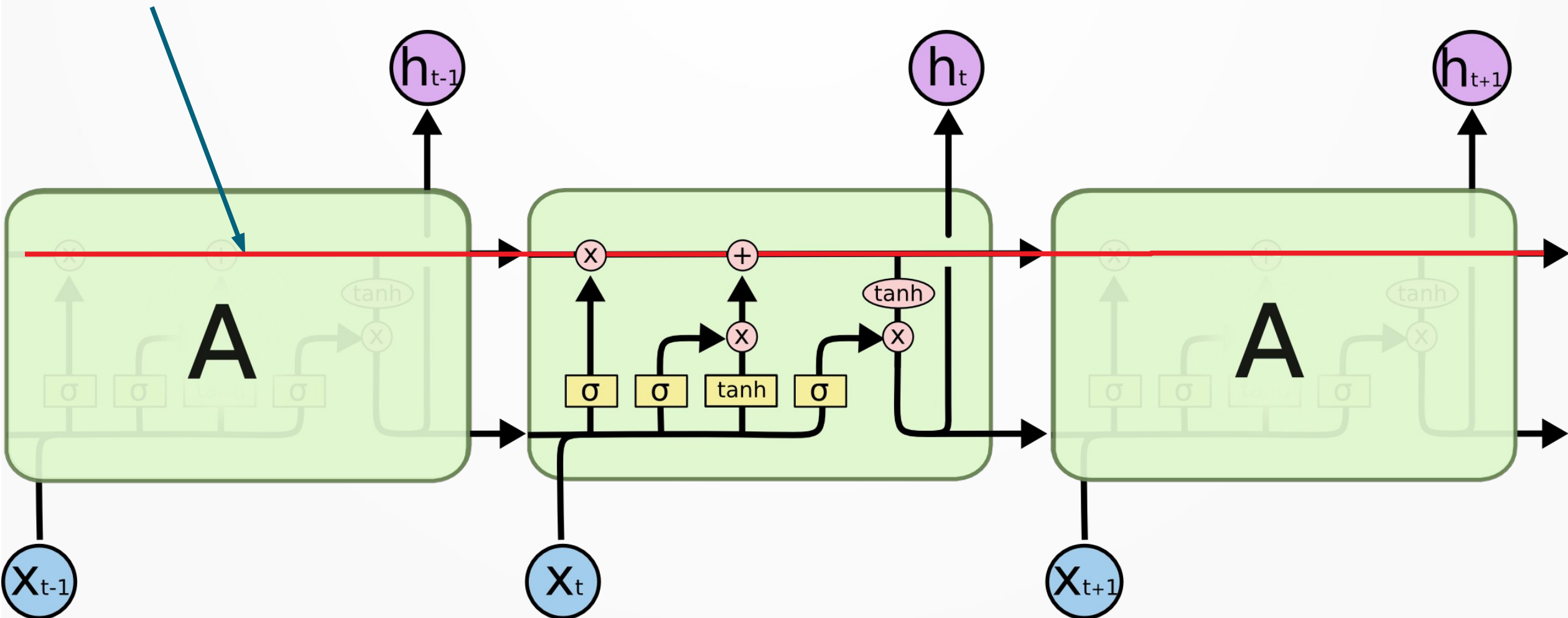
# Standard RNN





# LSTM

A LSTM cell can keep its internal state unchanged over many timesteps. This way a network can learn long term dependencies.



# further information

- Chris Olah: [Understanding LSTM Networks](#)
- Jürgen Schmidhuber: [Tutorial on LSTM Recurrent Networks](#)
- **LSTM: A Search Space Odyssey**  
Klaus Greff, Rupesh K. Srivastava, Jan Koutník, Bas R. Steunebrink, Jürgen Schmidhuber 2015
- **Speech Recognition with Deep Recurrent Neural Networks**  
Graves et al. 2013
- **Long Short-term Memory**  
Sepp Hochreiter, Jürgen Schmidhuber

**What does the network learn?**

### Cell sensitive to position in line:

The sole importance of the crossing of the Berezina lies in the fact that it plainly and indubitably proved the fallacy of all the plans for cutting off the enemy's retreat and the soundness of the only possible line of action--the one Kutuzov and the general mass of the army demanded--namely, simply to follow the enemy up. The French crowd fled at a continually increasing speed and all its energy was directed to reaching its goal. It fled like a wounded animal and it was impossible to block its path. This was shown not so much by the arrangements it made for crossing as by what took place at the bridges. When the bridges broke down, unarmed soldiers, people from Moscow and women with children who were with the French transport, all--carried on by vis inertiae--pressed forward into boats and into the ice-covered water and did not, surrender.

Cell that turns on inside quotes:

"You mean to imply that I have nothing to eat out of.... On the contrary, I can supply you with everything even if you want to give dinner parties," warmly replied Chichagov, who tried by every word he spoke to prove his own rectitude and therefore imagined Kutuzov to be animated by the same desire.

Kutuzov, shrugging his shoulders, replied with his subtle penetrating smile: "I meant merely to say what I said."



Cell that robustly activates inside if statements:

```
static int __dequeue_signal(struct sigpending *pending, sigset_t *mask,
                           siginfo_t *info)
{
    int sig = next_signal(pending, mask);
    if (sig) {
        if (current->notifier) {
            if (sigismember(current->notifier_mask, sig)) {
                if (!(current->notifier)(current->notifier_data)) {
                    clear_thread_flag(TIF_SIGPENDING);
                    return 0;
                }
            }
        }
        collect_signal(sig, pending, info);
    }
    return sig;
}
```

A large portion of cells are not easily interpretable. Here is a typical example:

```
/* Unpack a filter field's string representation from user-space
 * buffer. */
char *audit_unpack_string(void **bufp, size_t *remain, size_t len)
{
    char *str;
    if (!*bufp || (len == 0) || (len > *remain))
        return ERR_PTR(-EINVAL);
    /* Of the currently implemented string fields, PATH_MAX
     * defines the longest valid length.
     */
}
```

# What does the network learn?

- Visualizing the predictions and the “neuron” firings in the RNN
- Set the background color based on the neurons activation
- Literature tip:
  - <http://karpathy.github.io/2015/05/21/rnn-effectiveness>
- Paper:
  - Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016



# Other applications for RNNs

- Time Series Prediction
- Speech Recognition
  - Speech Recognition with Deep Recurrent Neural Networks Graves et al. 2013
- Drawing (Pictures, handwriting)
  - Generating Sequences With Recurrent Neural Networks Graves 2013
- Music Generation
  - Song From PI: A Musically Plausible Network for Pop Music Generation

# Summary

- RNNs are used to deal with sequential data
- They have applications in many other domains:
  - Speech Recognition
  - Time Series Prediction
  - Drawing (Pictures, Handwriteting)
- During Backpropagation gradients can explode or vanish
- Deep NLU is a hot topic of research (Bert, GTP-2)

**Write like  
Shakespeare**



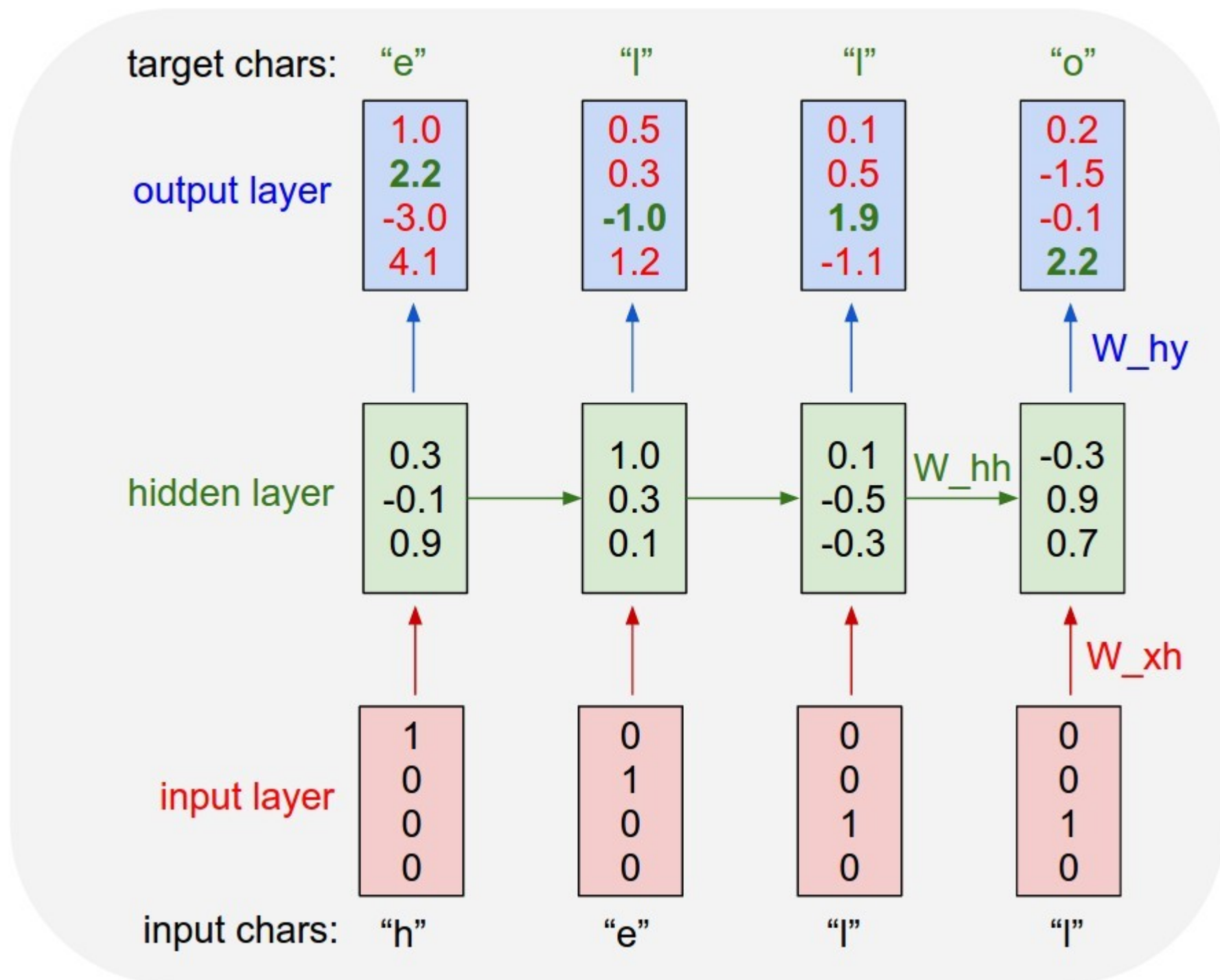


**Write like  
Shakespeare  
with character  
level RNNs**



# Tutorial

- We build a simple character level language model
- What you learn:
  - One Hot Encoding
  - Using LSTM and GRU
  - Character level prediction





# Questions

- How does the a change of training epochs effect the result?
- What happens when you add more layers?
- How does the result change when you replace LSTM with GRUs?
- How does this network work for other texts and languages?
- Do these text make any sense? Explain why.
- Send me interesting results: [oliver.guhr@htw-dresden.de](mailto:oliver.guhr@htw-dresden.de)

