

# Supervised Learning: Deep(er) learning

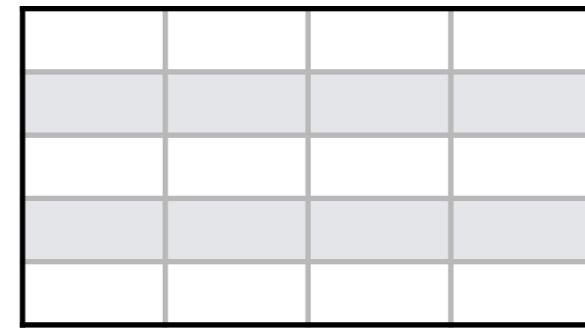
Jean Feng & Ali Shojaie

Aug 19-21, 2024  
Summer Institute in Statistics for Big Data  
University of Washington

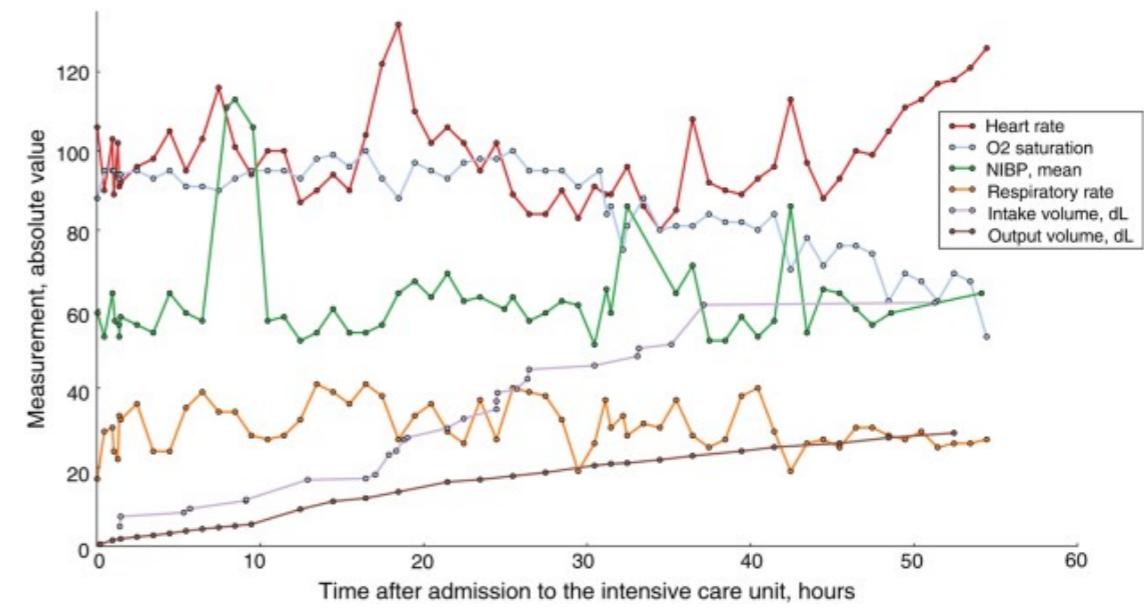
# Sequence data

- Up to now, we've assumed all inputs to a machine learning model have the same exact size
  - Images
  - Tabular data
- What happens if we are trying to predict labels using sequential data, e.g.
  - Text data
  - Time series data
  - Genetic sequences
  - Video analysis
  - ...

*Tabular data:*

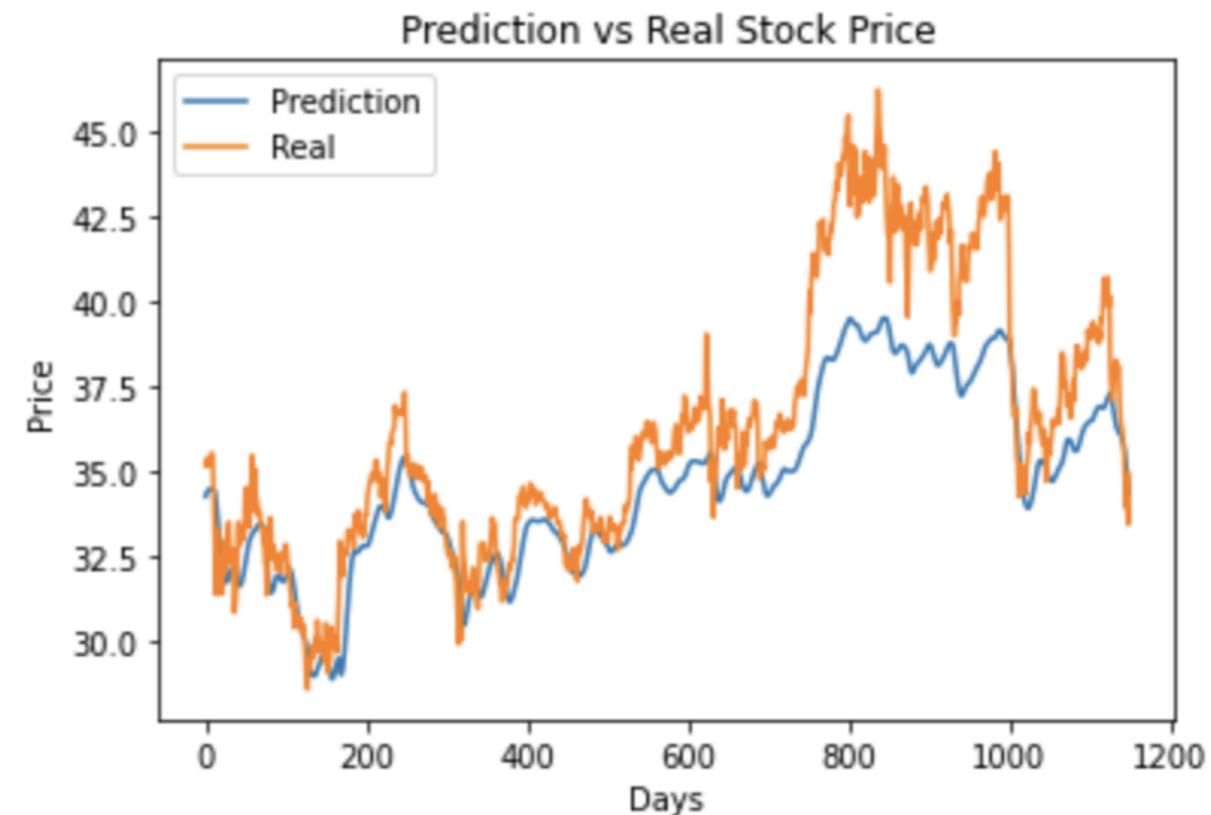


*Patient data from the ICU:*



# Example tasks for sequence data

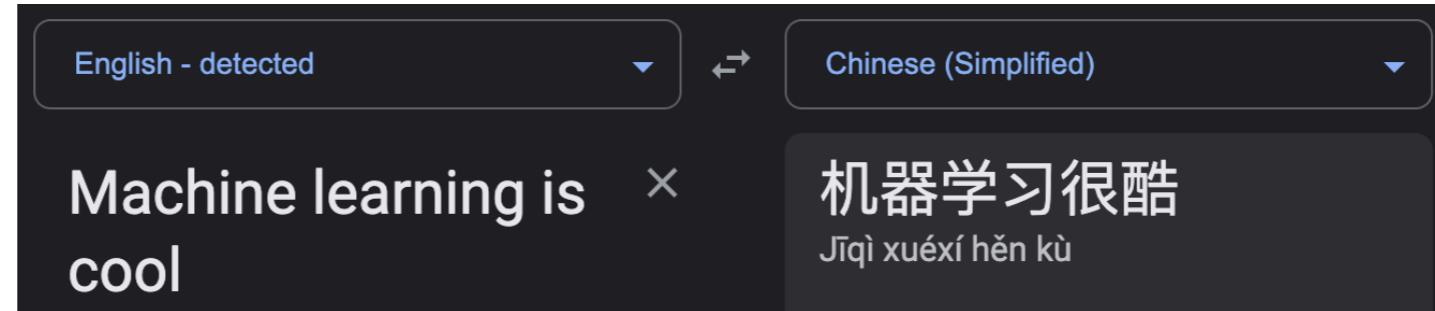
- **Auto-complete:**
  - *Text:* The patient said \_\_\_\_\_
- *Time series data:*



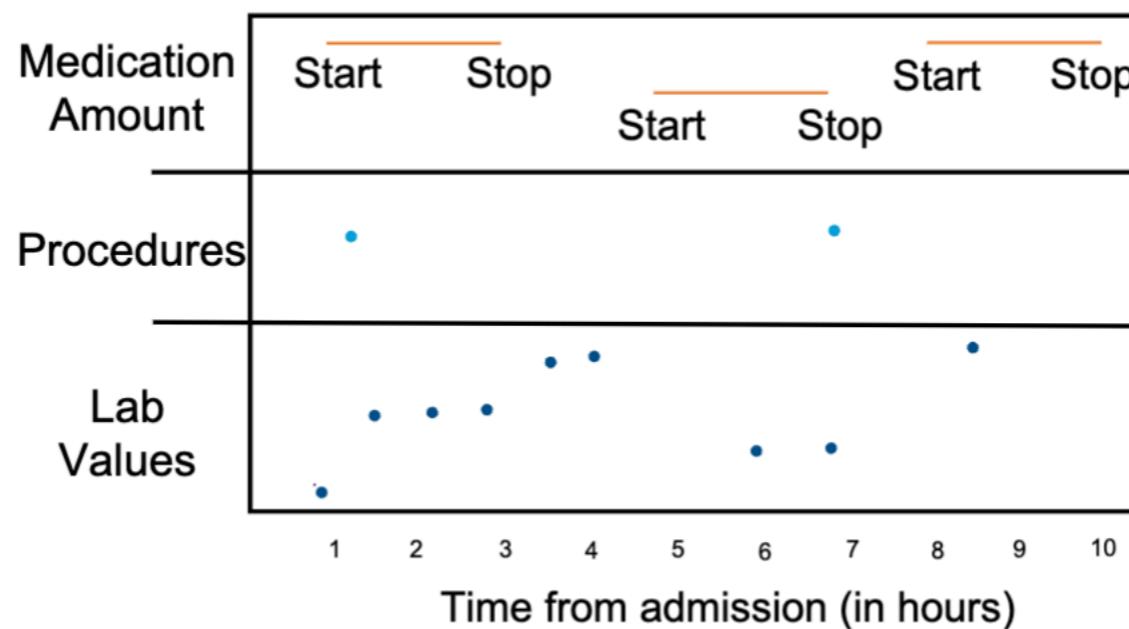
*Do this well and you'll be rich! ;)*

# Example tasks for sequence data

- Translation:



- Classification:



*Readmission  
within 30-days?*

# Example tasks for sequence data

Knowing When to Look: Adaptive Attention via  
A Visual Sentinel for Image Captioning

Jiasen Lu<sup>2\*</sup>, Caiming Xiong<sup>1†</sup>, Devi Parikh<sup>3</sup>, Richard Socher<sup>1</sup>

<sup>1</sup>Salesforce Research, <sup>2</sup>Virginia Tech, <sup>3</sup>Georgia Institute of Technology

jiasenlu@vt.edu, parikh@gatech.edu, {cxiong, rsocher}@salesforce.com

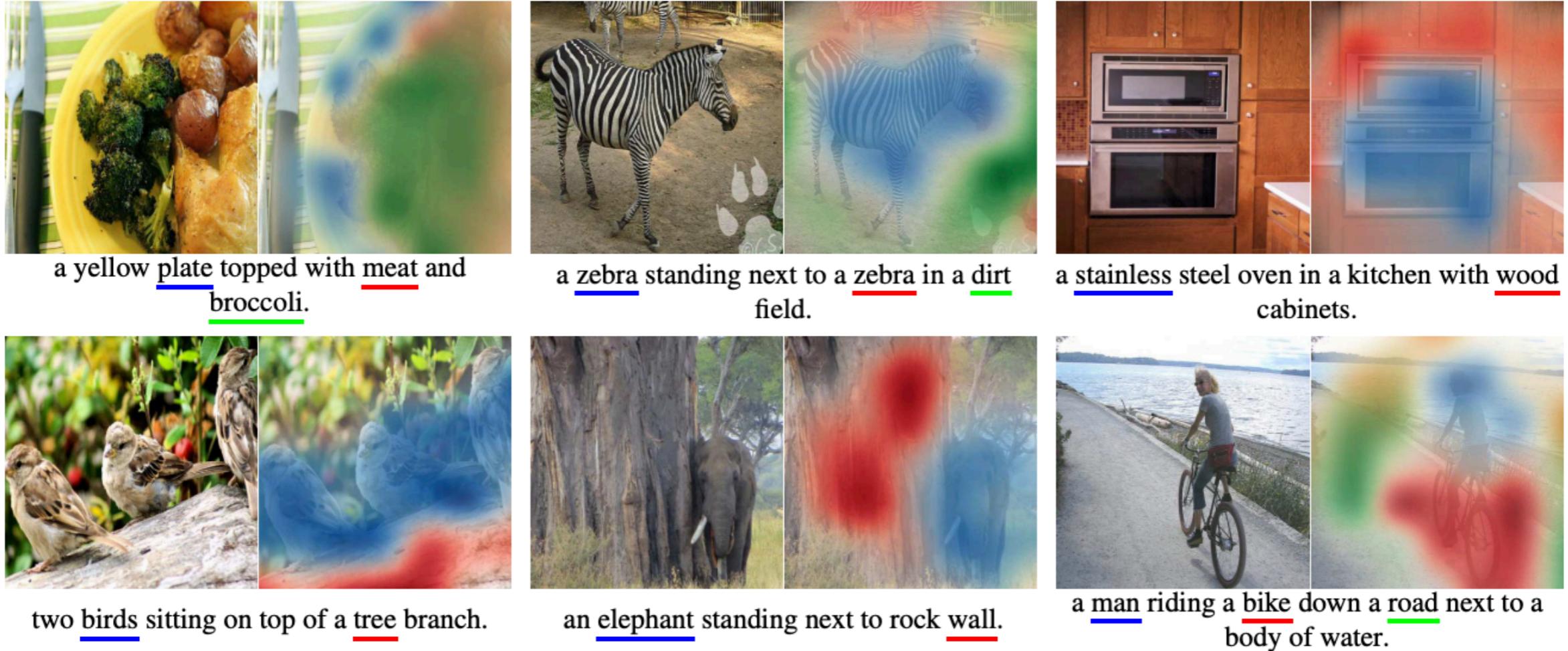


Figure 4: Visualization of generated captions and image attention maps on the COCO dataset. Different colors show a correspondence between attended regions and underlined words. First 2 rows are success cases, last rows are failure examples. Best viewed in color.

Lu et al, CVPR 2017

# How do humans process text data?

*Does this refer to a brain network? A gene network?*



+



**Deep learning** is the subset of **machine learning** methods based on **neural networks** with **representation learning**. The adjective "deep" refers to the use of multiple layers in the **network**.

Methods used can be either **supervised**, **semi-supervised** or **unsupervised**.<sup>[2]</sup>

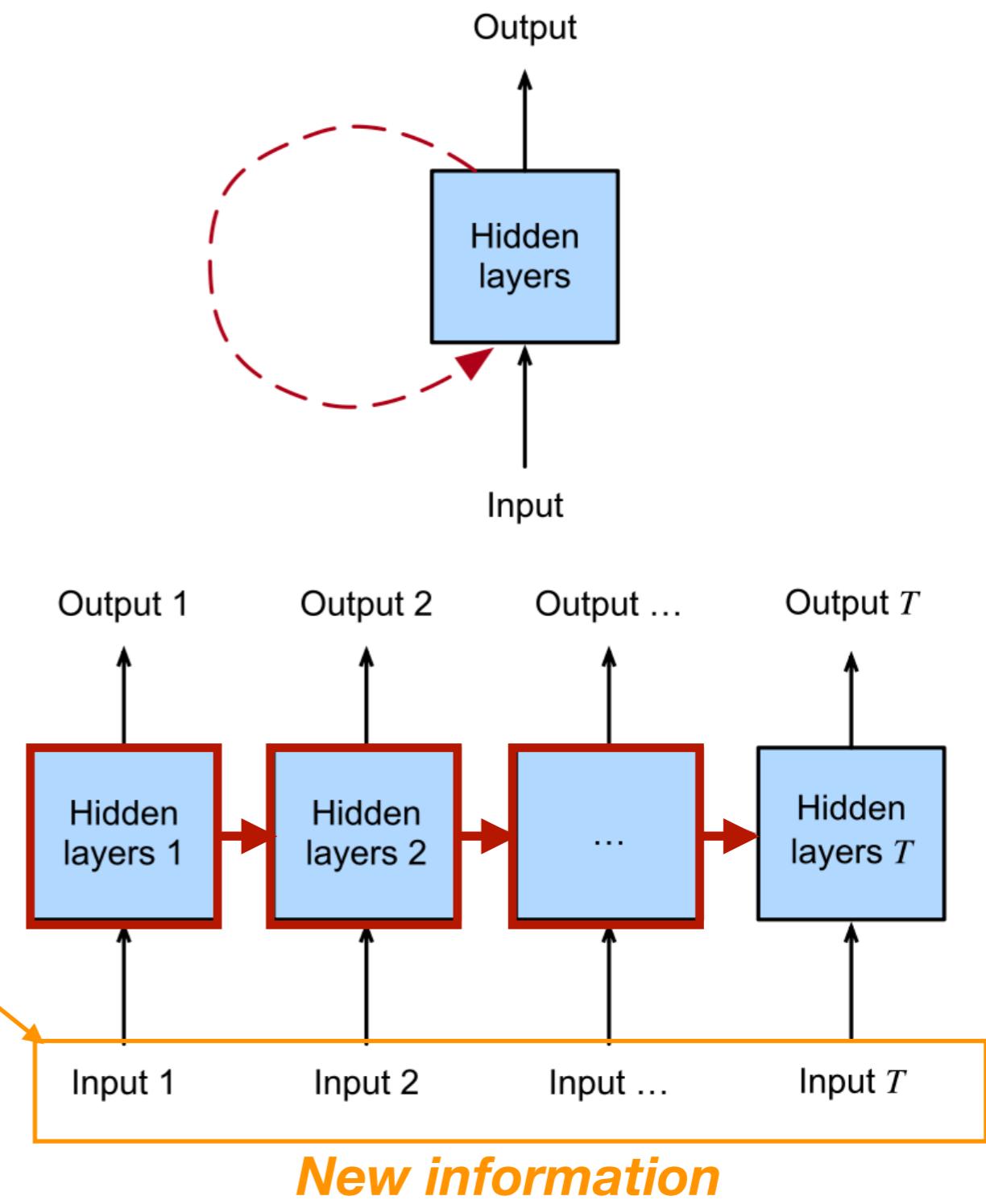
Deep-learning architectures such as **deep neural networks**, **deep belief networks**, **recurrent neural networks**, **convolutional neural networks** and **transformers** have been applied to fields including computer vision, speech recognition, natural language processing, machine translation, bioinformatics, drug design, medical image analysis, climate science, material inspection and board game programs, where they have produced results comparable to and in some cases surpassing human expert performance.<sup>[3][4][5]</sup>

Early forms of neural networks were inspired by information processing and distributed communication nodes in **biological systems**, in particular the **human brain**. However, current neural networks do not intend to model the brain function of organisms, and are generally seen as low quality models for that purpose.<sup>[6]</sup>

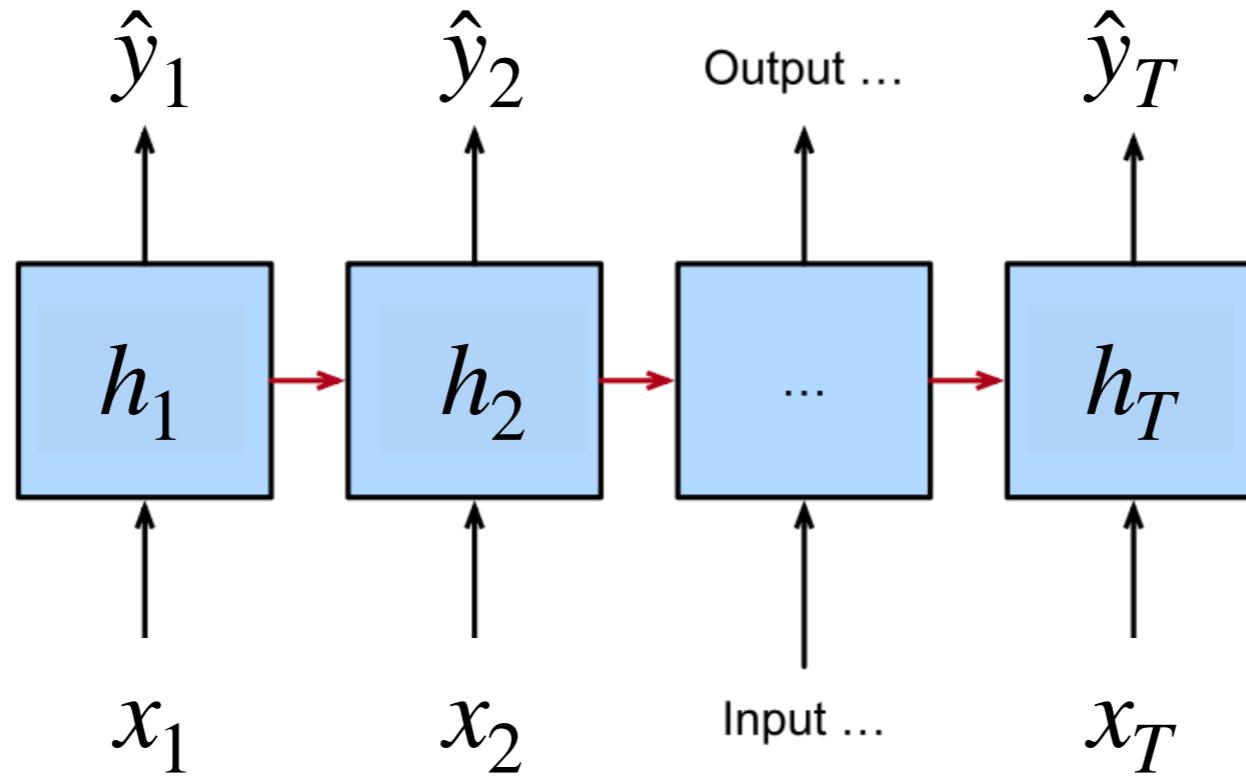
# Recurrent neural networks

- **Key idea:** Define an NN that will be applied in the same manner to each position in the sequence. To make a prediction at position  $t$ , the inputs to the RNN are:

- The input at position  $t$
- (Intermediate) Outputs at position  $t - 1$   
“Memory”



# Any ideas how we make this into math?

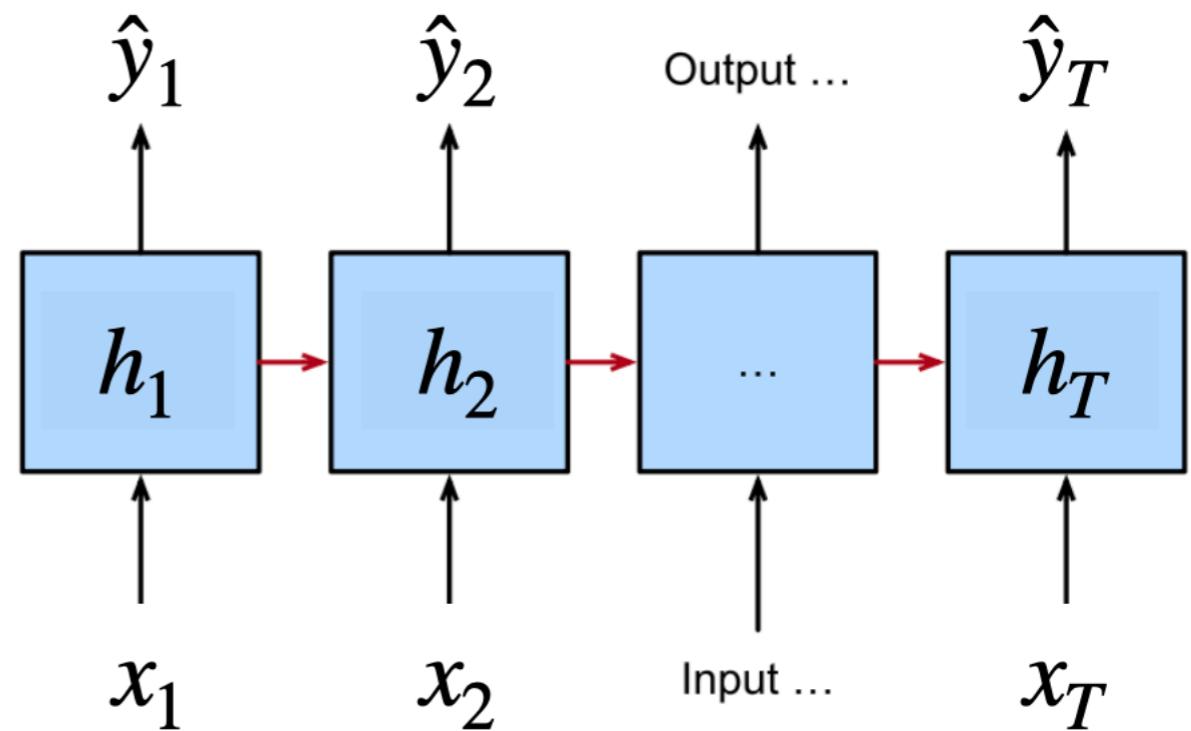


- $x_t$  = Input at position  $t$
- $\hat{y}_t$  = Prediction at position  $t$
- $h_t$  = Hidden layer at position  $t$ , encodes “memory”

# Recurrent neural networks

## Notation:

- $x_t$  = Input at position  $t$
- $\hat{y}_t$  = Prediction at position  $t$
- $h_t$  = Hidden layer at position  $t$ , encodes “memory”



## Example of a simple RNN:

$$h_t = \sigma_1(Uh_{t-1} + Wx_t)$$

$$\hat{y}_t = \sigma_2(Vh_t)$$

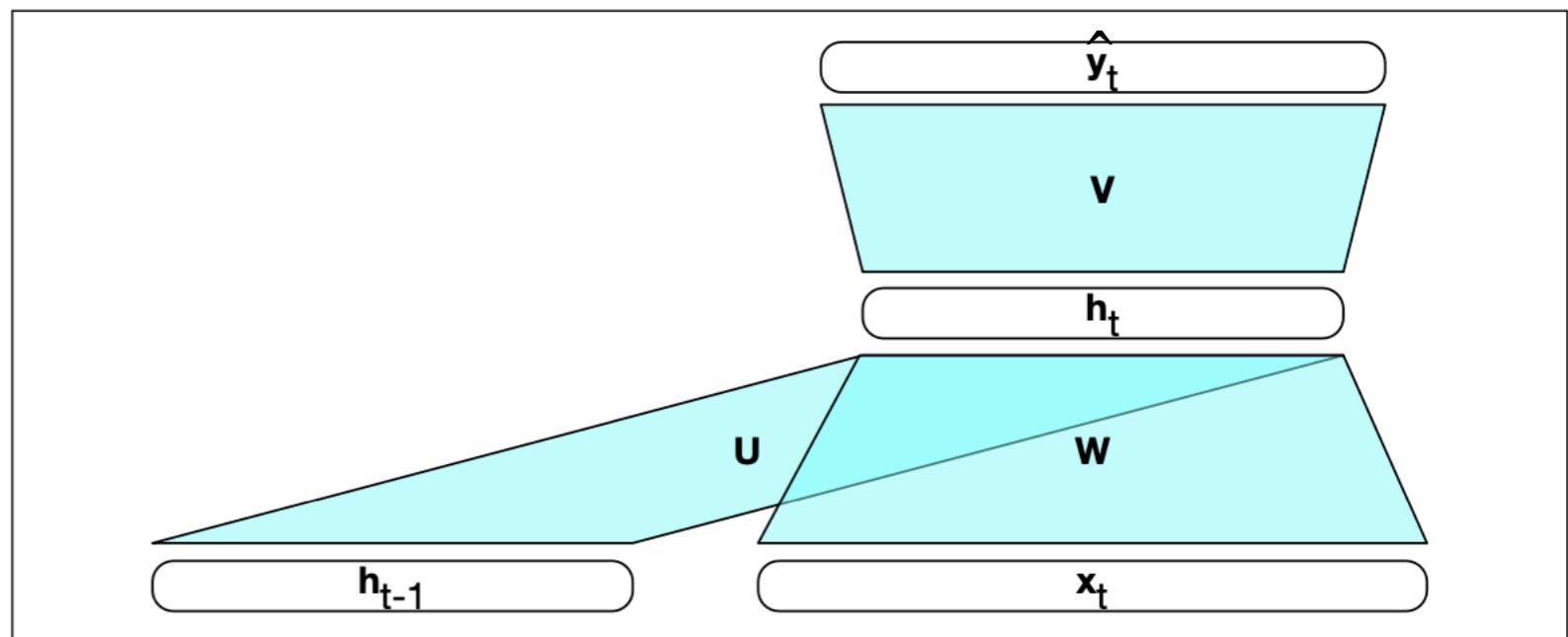


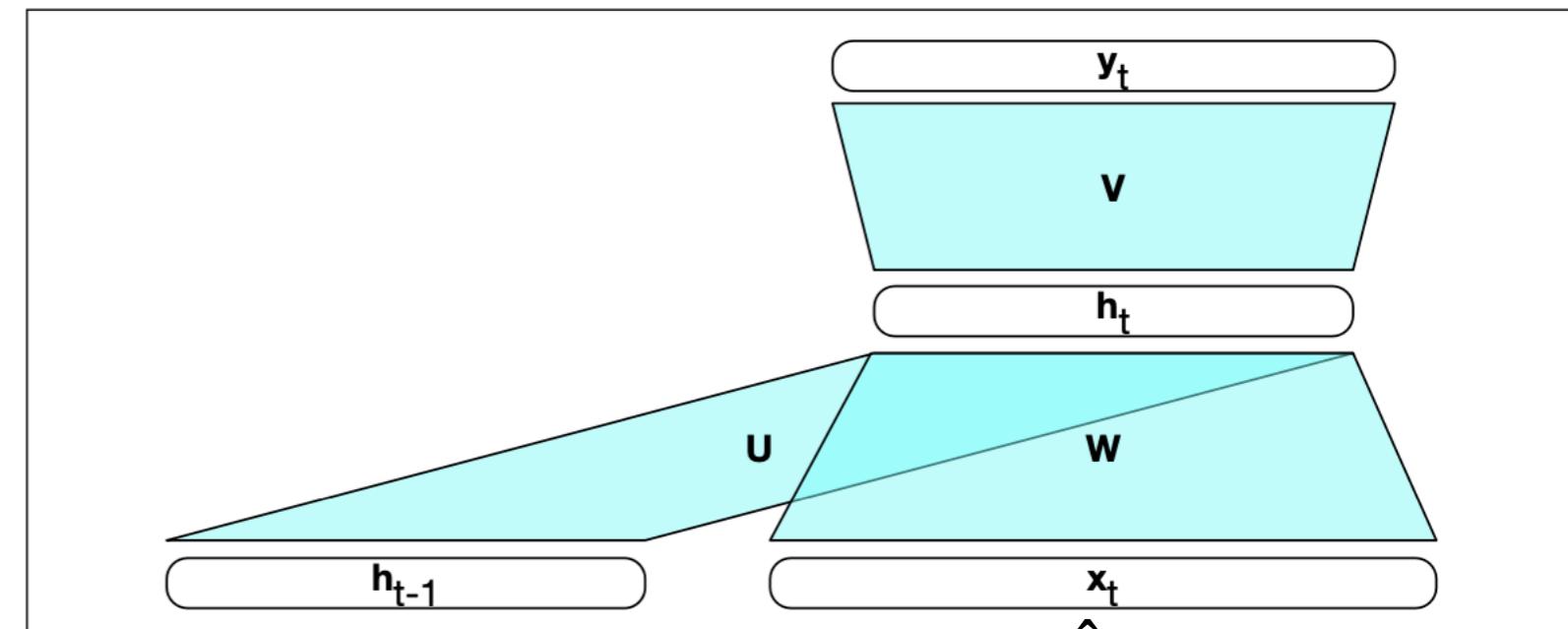
Figure 9.2 Simple recurrent neural network illustrated as a feedforward network.

# Recurrent neural networks

**Example of a simple RNN:**

$$h_t = \sigma_1(Uh_{t-1} + Wx_t)$$

$$\hat{y}_t = \sigma_2(Vh_t)$$



**Figure 9.2** Simple recurrent neural network illustrated as a feedforward network.

*What are the parameters that need to be learned in this RNN?*

# RNNs for text

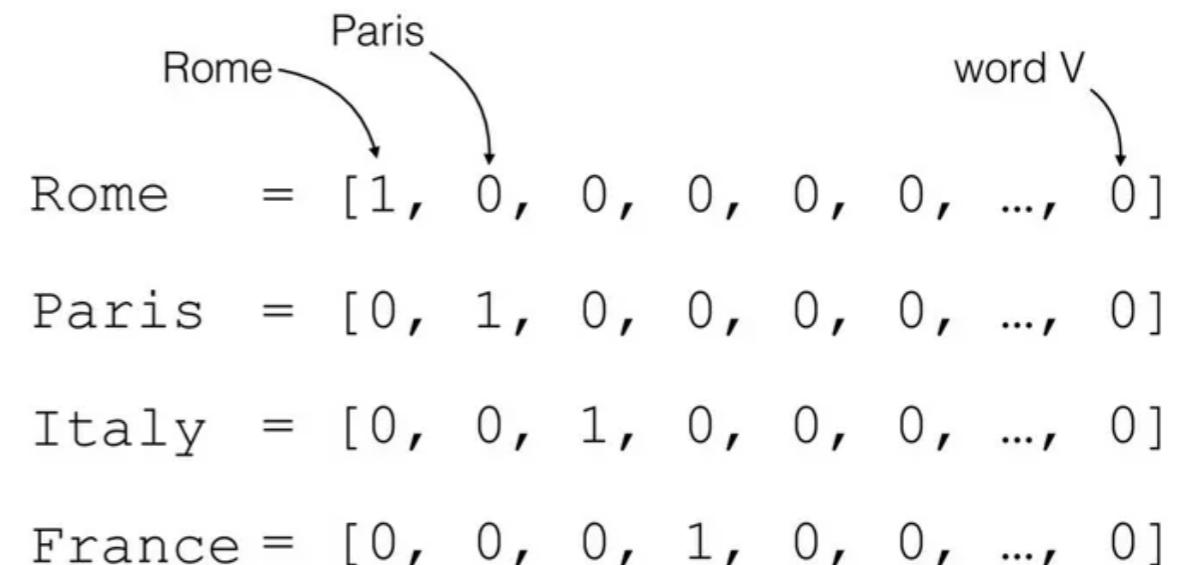
- Ok... but how do you encode a word as an input vector?
  - How do we represent words to an RNN (i.e. input)?
  - How do we words generated by an RNN (i.e. output)?

## *Example dataset*

If I smelled the scent of hand sanitizers today on someone in the past, I would think they were so i...

How #COVID19 Will Change Work in General (and recruiting, specifically) via/ @ProactiveTalent #Recru...

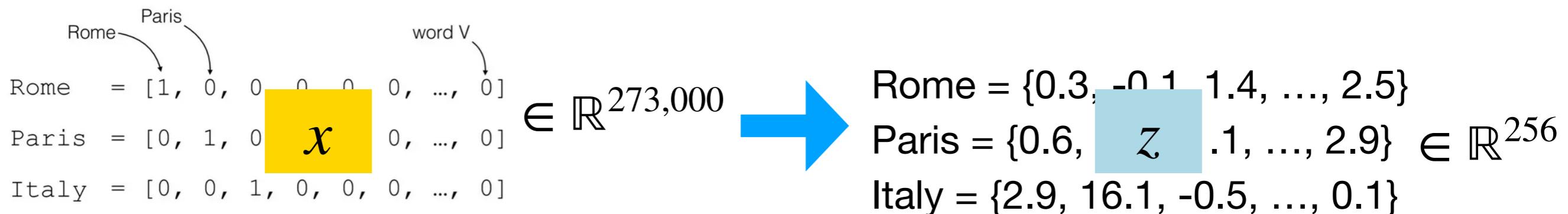
#coronavirus  
#covid19 deaths continue to rise. It's almost as bad as it ever was. Politicians and ...



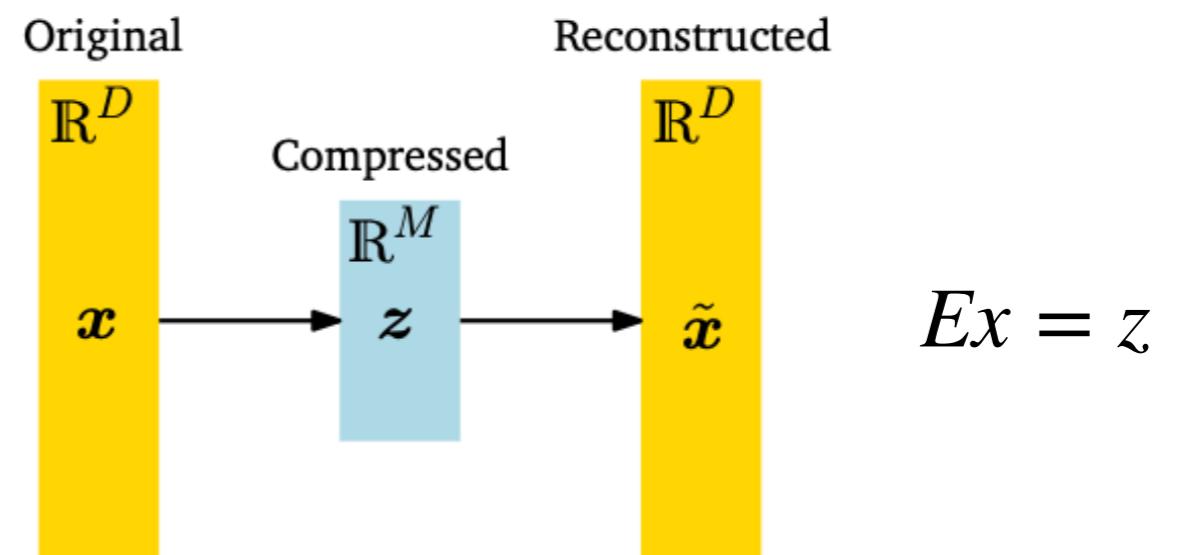
# RNNs for text

However, there are A LOT of words, so a one-hot encoding would be very high-dimensional.

What can we try to reduce the dimension of a high-dimensional dataset?



*How about PCA?*



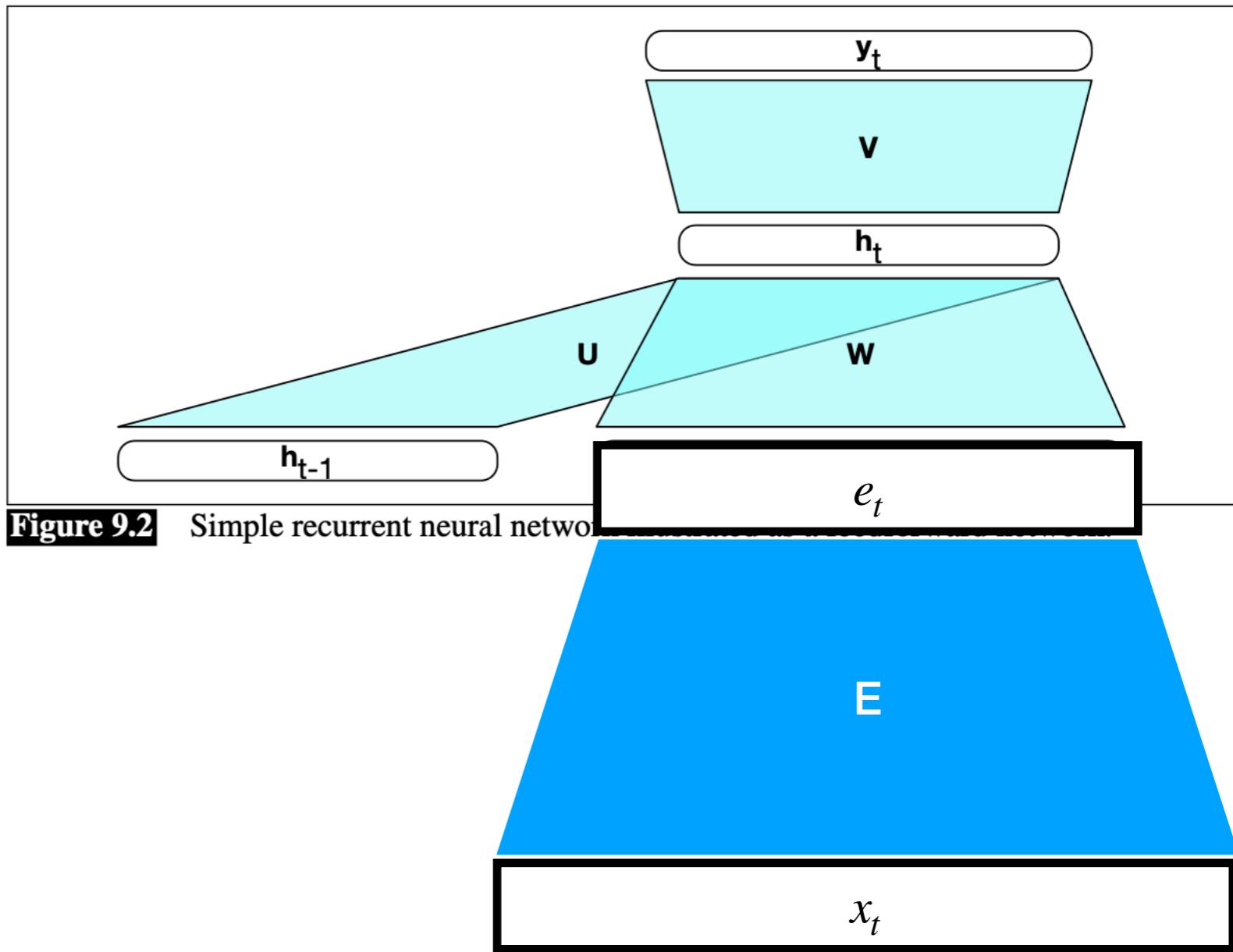
# RNNs for predicting text

**Example of a simple RNN for text:**

$$e_t = Ex_t$$

$$h_t = \text{softmax}(Uh_{t-1} + We_t)$$

$$\hat{p}_t = \text{softmax}(Vh_t)$$



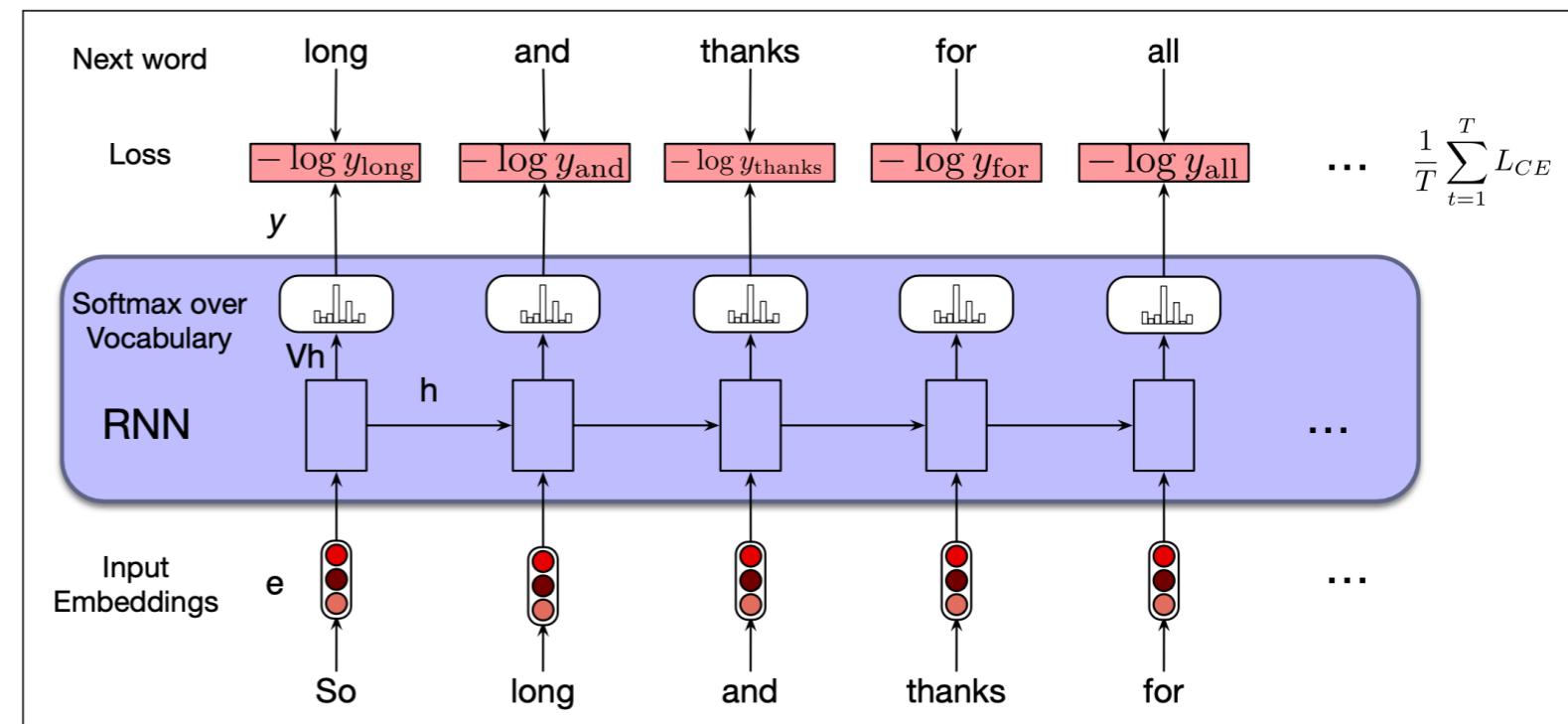
# RNNs for predicting text

**Example of a simple RNN for text:**

$$e_t = Ex_t$$

$$h_t = \text{softmax}(Uh_{t-1} + We_t)$$

$$\hat{p}_t = \text{softmax}(Vh_t)$$



**Figure 9.6** Training RNNs as language models.

# Word2vec

## Efficient Estimation of Word Representations in Vector Space

**Tomas Mikolov**  
Google Inc., Mountain View, CA  
tmikolov@google.com

**Greg Corrado**  
Google Inc., Mountain View, CA  
gcorrado@google.com

**Kai Chen**  
Google Inc., Mountain View, CA  
kaichen@google.com

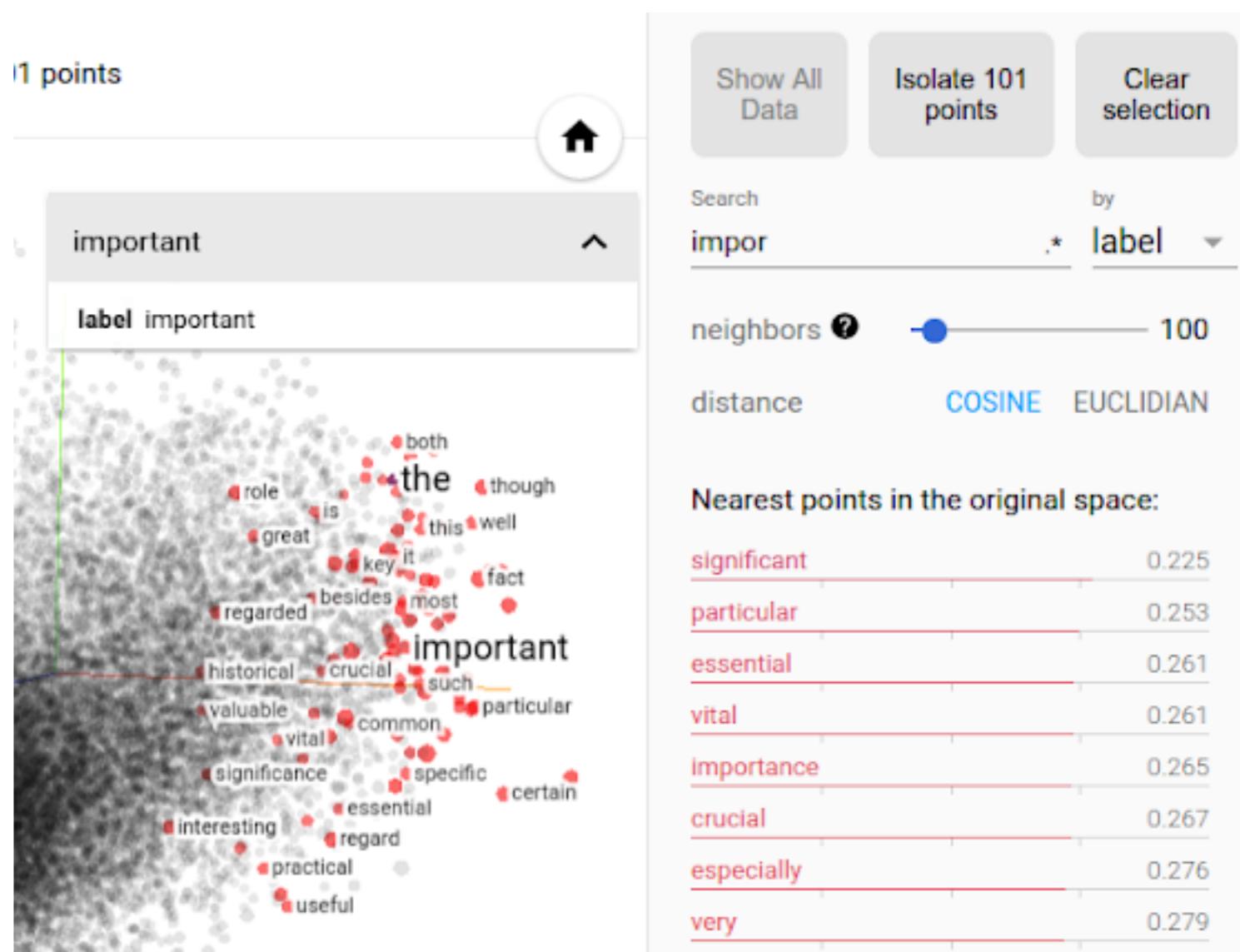
**Jeffrey Dean**  
Google Inc., Mountain View, CA  
jeff@google.com

### Abstract

We propose two novel model architectures for computing continuous vector representations of words from very large data sets. The quality of these representations is measured in a word similarity task, and the results are compared to the previously best performing techniques based on different types of neural networks. We observe large improvements in accuracy at much lower computational cost, i.e. it takes less than a day to learn high quality word vectors from a 1.6 billion words data set. Furthermore, we show that these vectors provide state-of-the-art performance on our test set for measuring syntactic and semantic word similarities.

**Training data:** internal Google dataset with 1+ billion words.

**Vocabulary:** 692K words, after discarding all words that occurred less than 5 times in the training data



<https://projector.tensorflow.org/>

# Word2vec

## Linguistic Regularities in Continuous Space Word Representations

Tomas Mikolov\*, Wen-tau Yih, Geoffrey Zweig

Microsoft Research  
Redmond, WA 98052

### Abstract

Continuous space language models have recently demonstrated outstanding results across a variety of tasks. In this paper, we examine the vector-space word representations that are implicitly learned by the input-layer weights. We find that these representations are surprisingly good at capturing syntactic and semantic regularities in language, and that each relationship is characterized by a relation-specific vector offset. This allows vector-oriented reasoning based on the offsets between words. For example, the male/female relationship is automatically learned, and with the induced vector representations, “King - Man + Woman” results in a vector very close to “Queen.” We demonstrate that the word vectors capture syntactic regularities by means of syntactic analogy questions (provided with this paper), and are able to correctly answer almost 40% of the questions. We demonstrate that the word vectors capture semantic regularities by using the vector offset method to answer SemEval-2012 Task 2 questions. Remarkably, this method outperforms the best previous systems.

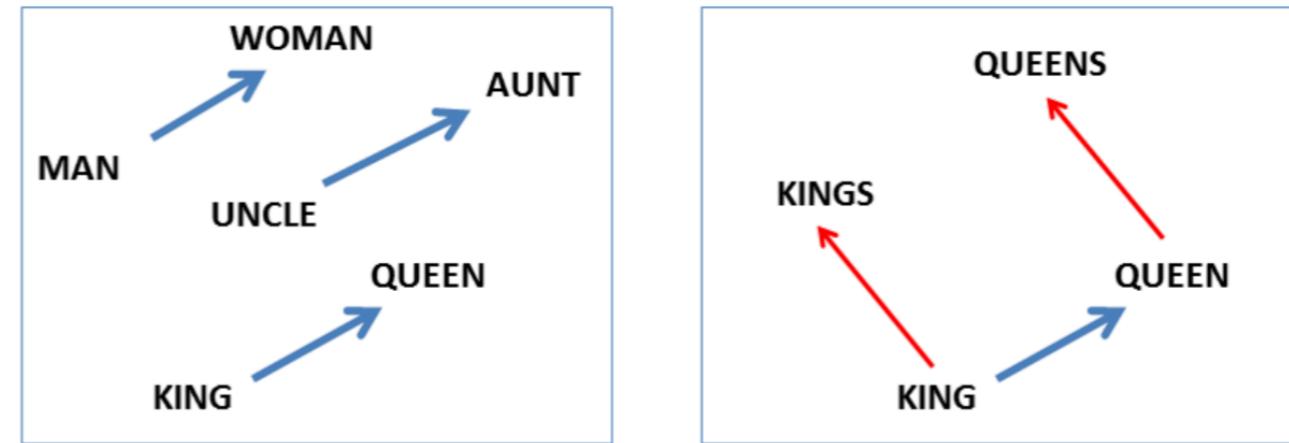


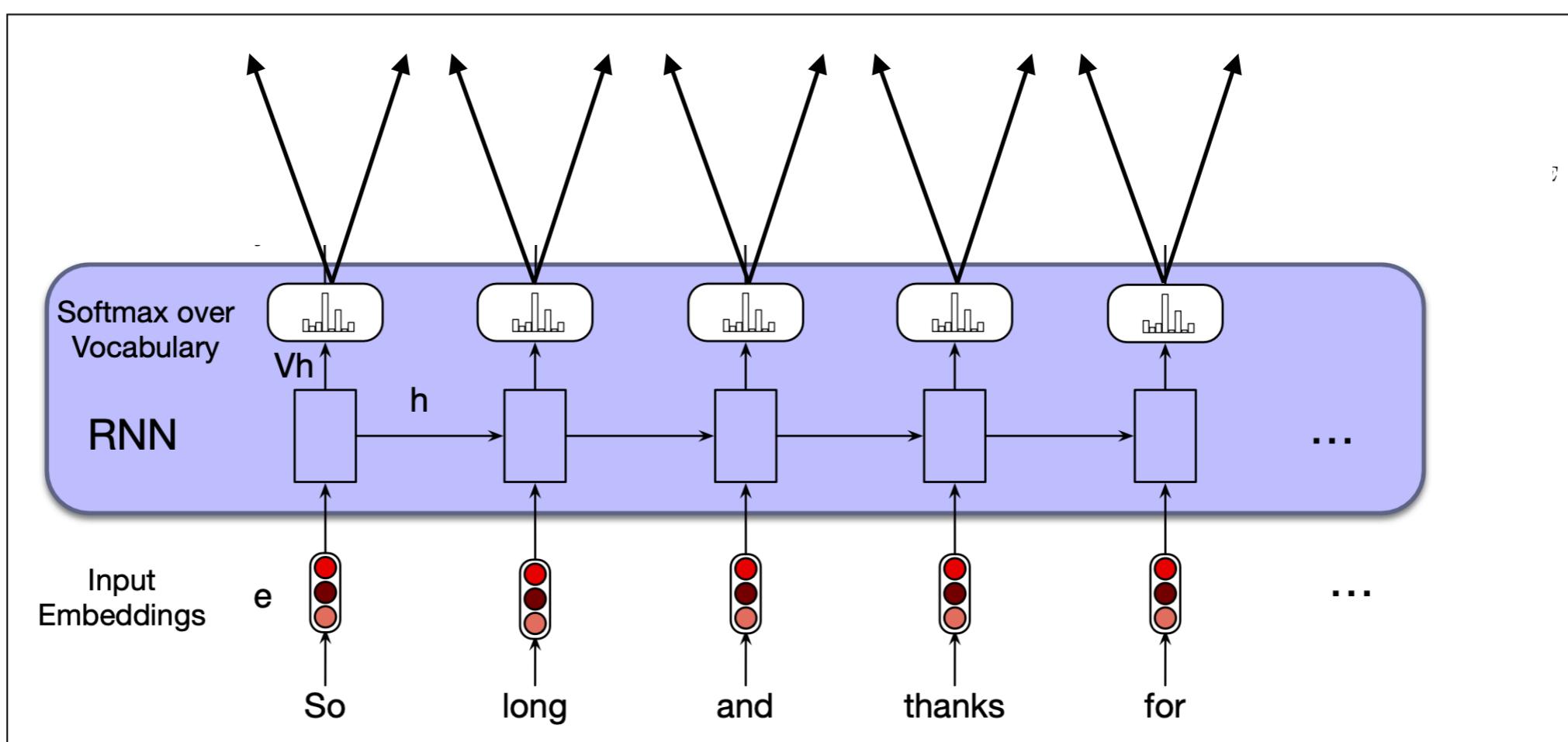
Figure 2: Left panel shows vector offsets for three word pairs illustrating the gender relation. Right panel shows a different projection, and the singular/plural relation for two words. In high-dimensional space, multiple relations can be embedded for a single word.

# Word2vec

- Now that we have embeddings for words, we can use that for lots of other tasks!
  - **Text classification:** Encode a sentence/phrase by the average word embedding
  - **Initialize embeddings for other models**
  - **Knowledge discovery:** By analyzing embeddings, you can discover new knowldge, e.g. **relationships between genes and diseases**

# Your turn to make a NN!

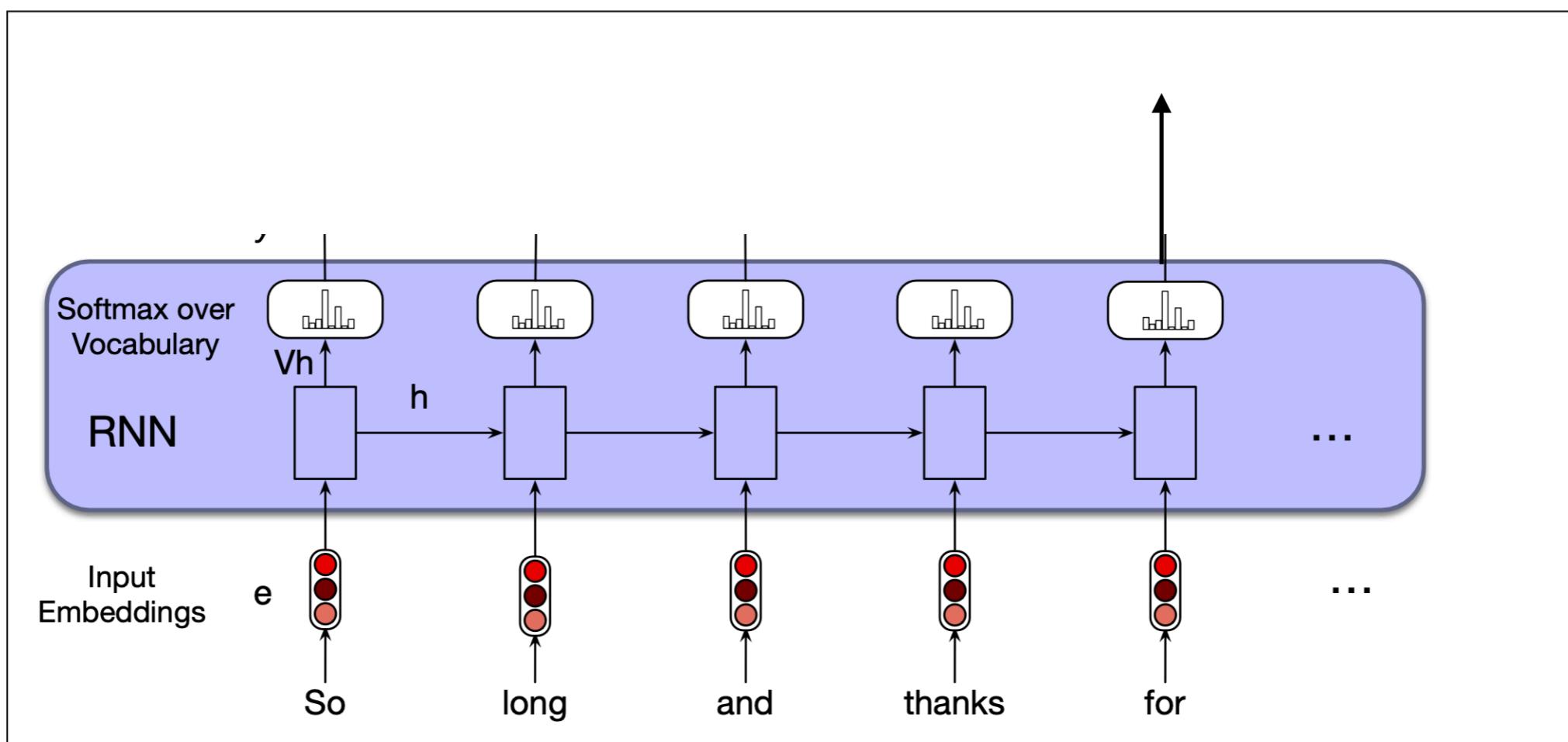
- Q: We've talked about how to create an RNN to predict the next word. How would you modify the RNN from before to predict the next two words?



**Figure 9.6** Training RNNs as language models.

# Your turn to make a NN!

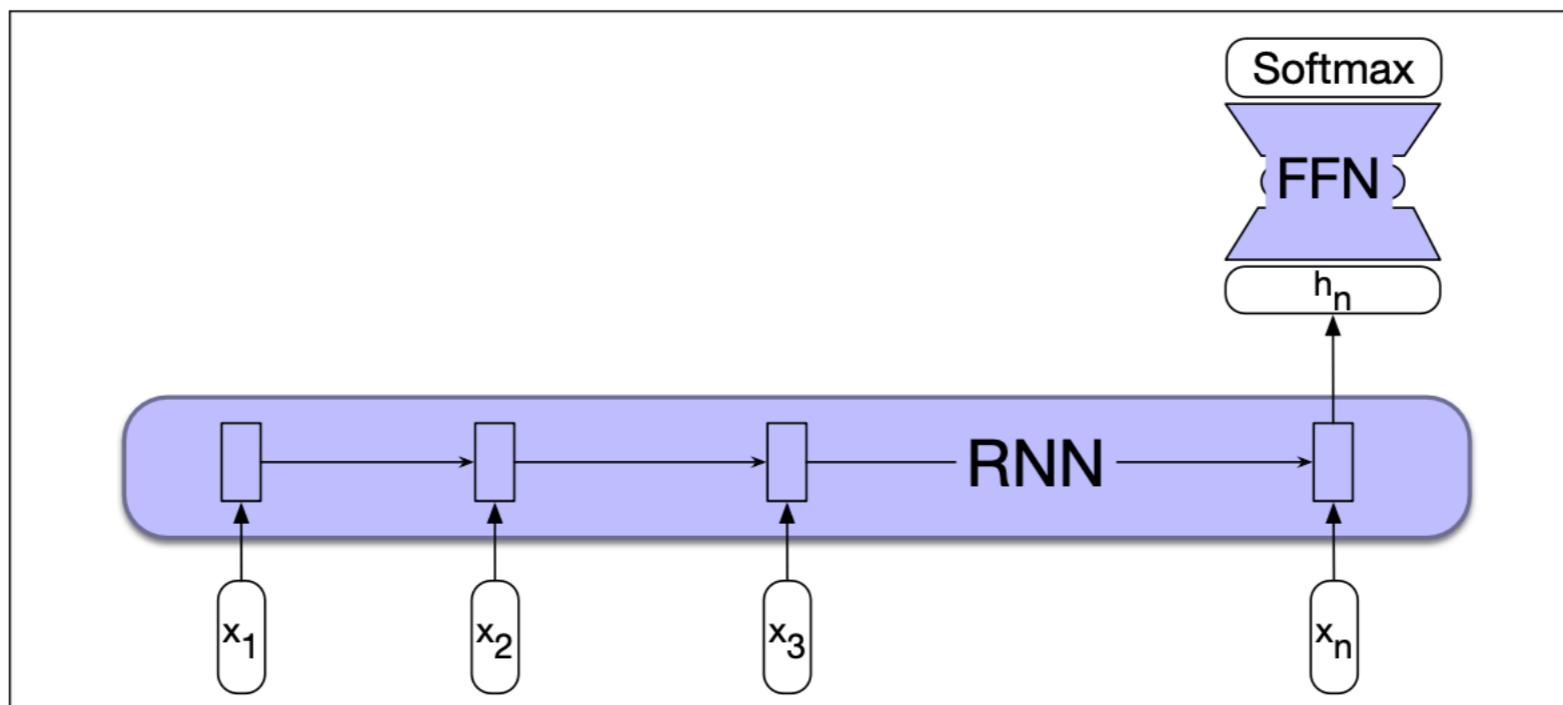
- Q: How would you construct an RNN for classifying sentences, e.g. outputs a single probability for a single sentence? How would you modify the RNN from before?



**Figure 9.6** Training RNNs as language models.

# Your turn to make a NN!

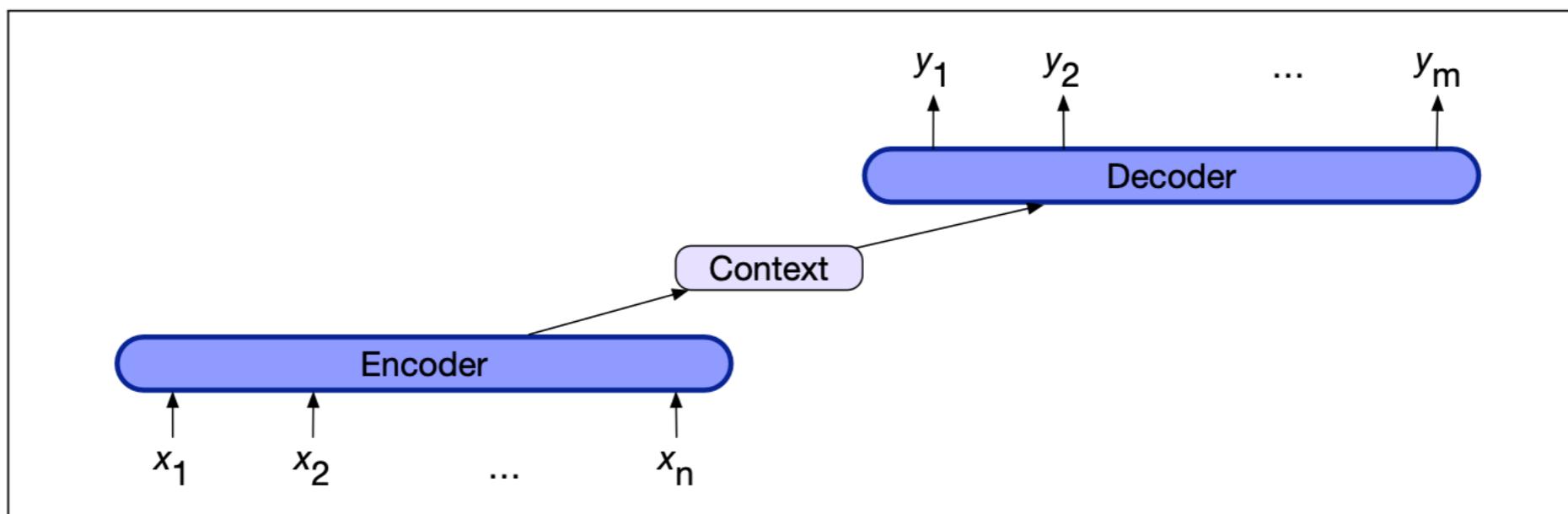
- Q: How would you construct an RNN for classifying sentences, e.g. outputs a single probability for a single sentence? How would you modify the RNN from before?



**Figure 9.8** Sequence classification using a simple RNN combined with a feedforward network. The final hidden state from the RNN is used as the input to a feedforward network that performs the classification.

# Your turn to make a NN!

- Q: What if we want to translate a sentence from English to a different language? Brainstorm an RNN architecture for conducting this translation task.



**Figure 9.16** The encoder-decoder architecture. The context is a function of the hidden representations of the input, and may be used by the decoder in a variety of ways.

# Outline

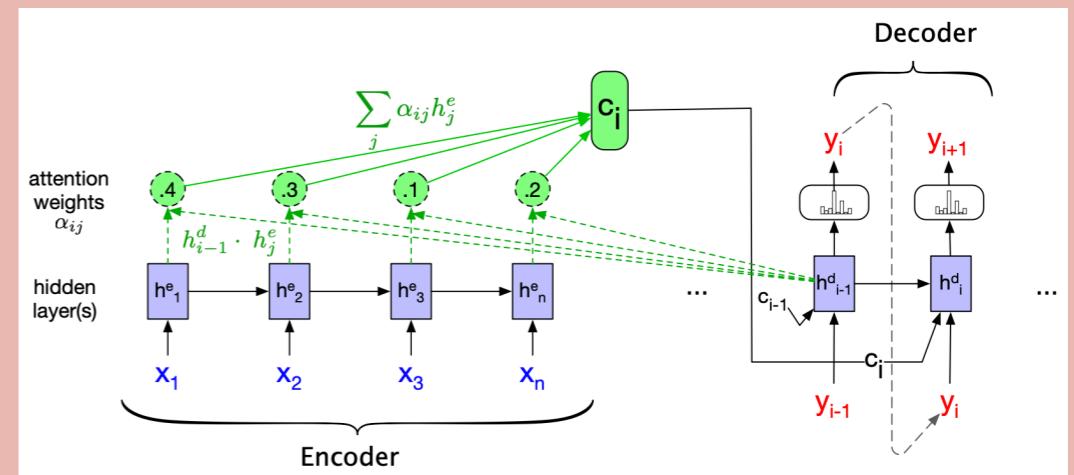
- Recurrent Neural Networks
- **Transformers**

# Inductive biases in NN

- Until very recently, the assumption was that we need inductive biases in the NN to help it generalize.
- **Convolutional neural networks:** Encodes translational invariance
- **Recurrent neural networks:** Also encodes translational invariance, but with memory

Problems with these models:

- Computation is slow



- Learning dependencies between distant positions is difficult

# What if we threw it all out the window?

## Attention Is All You Need

**Ashish Vaswani\***  
Google Brain  
avaswani@google.com

**Noam Shazeer\***  
Google Brain  
noam@google.com

**Niki Parmar\***  
Google Research  
nikip@google.com

**Jakob Uszkoreit\***  
Google Research  
usz@google.com

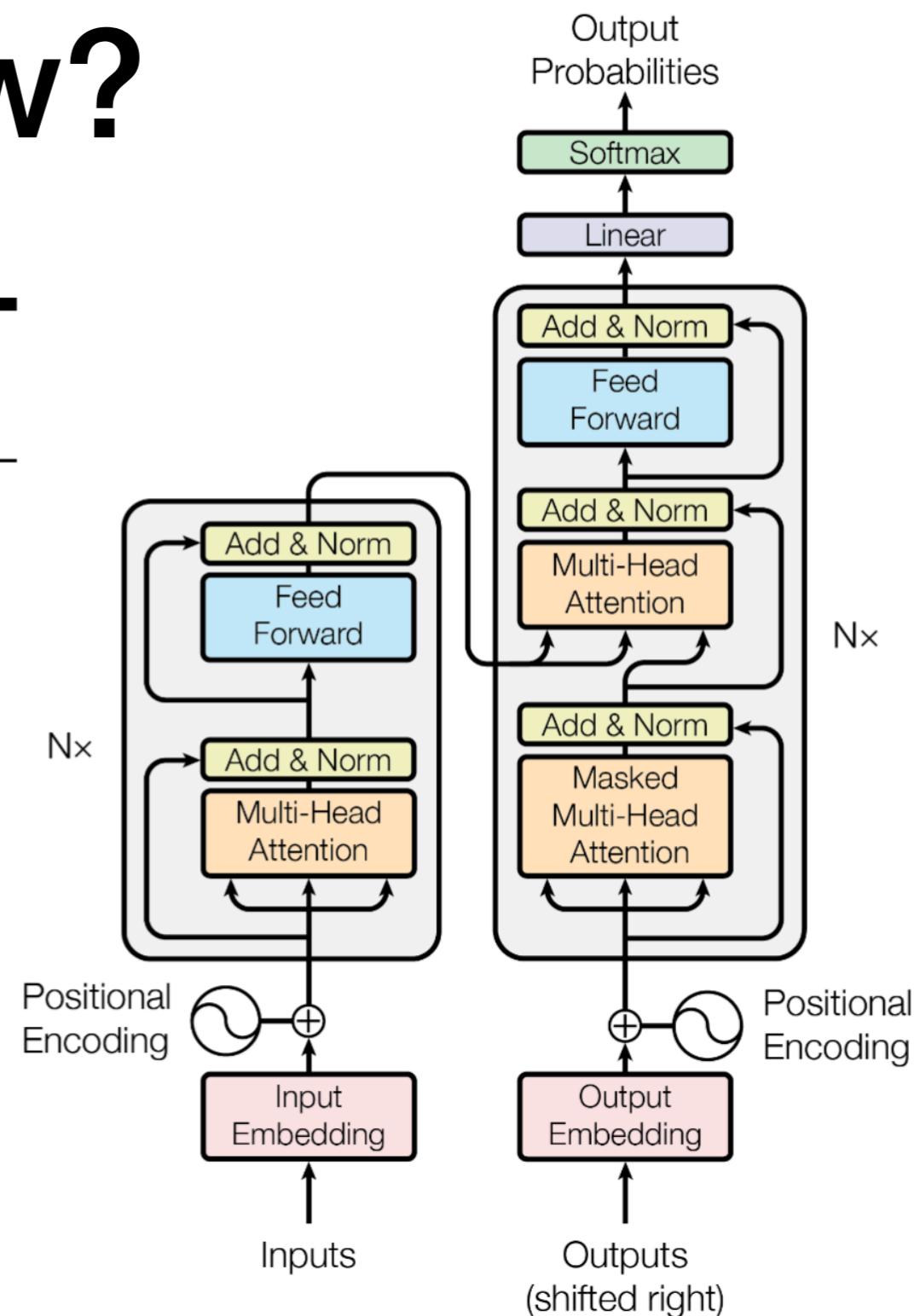
**Llion Jones\***  
Google Research  
llion@google.com

**Aidan N. Gomez\* †**  
University of Toronto  
aidan@cs.toronto.edu

**Łukasz Kaiser\***  
Google Brain  
lukaszkaiser@google.com

**Illia Polosukhin\* ‡**  
illia.polosukhin@gmail.com

*Essentially zero inductive biases*



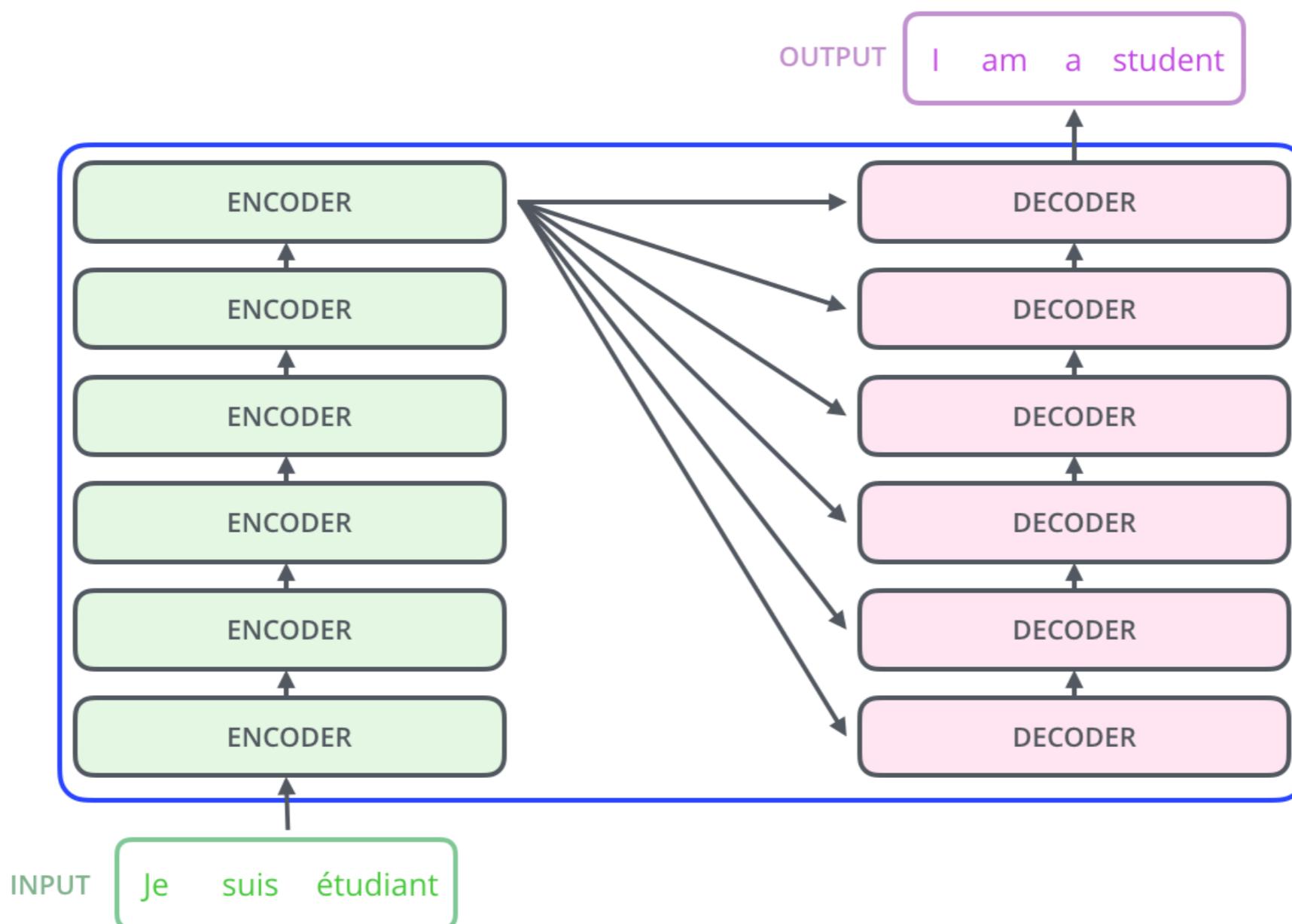
# Shoot, it works well.

We trained on the standard WMT 2014 English-German dataset consisting of about 4.5 million sentence pairs. Sentences were encoded using byte-pair encoding [3], which has a shared source-target vocabulary of about 37000 tokens. For English-French, we used the significantly larger WMT 2014 English-French dataset consisting of 36M sentences and split tokens into a 32000 word-piece vocabulary [31]. Sentence pairs were batched together by approximate sequence length. Each training batch contained a set of sentence pairs containing approximately 25000 source tokens and 25000 target tokens.

Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

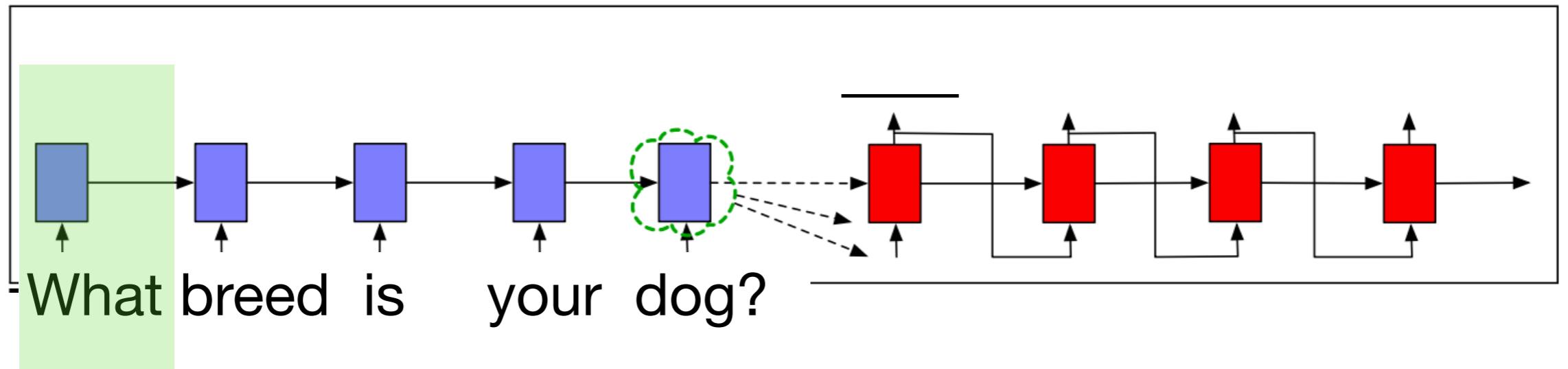
Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [15]	23.75			
Deep-Att + PosUnk [32]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [31]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [8]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [26]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [32]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [31]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [8]	26.36	<b>41.29</b>	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1		<b><math>3.3 \cdot 10^{18}</math></b>
Transformer (big)	<b>28.4</b>	<b>41.0</b>		$2.3 \cdot 10^{19}$

# Understanding transformers



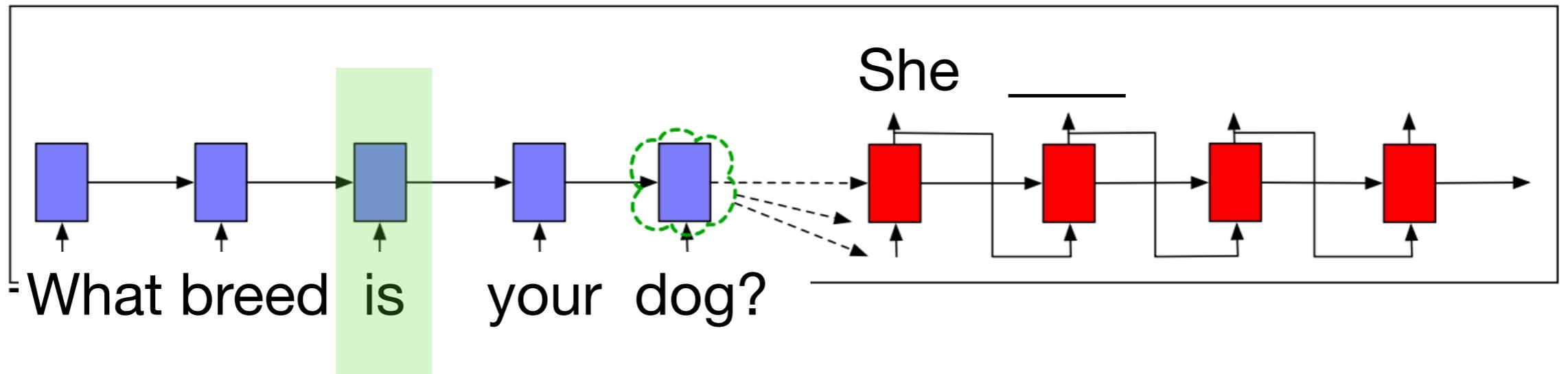
<https://jalammar.github.io/illustrated-transformer/>

# Attention: Which inputs matter?



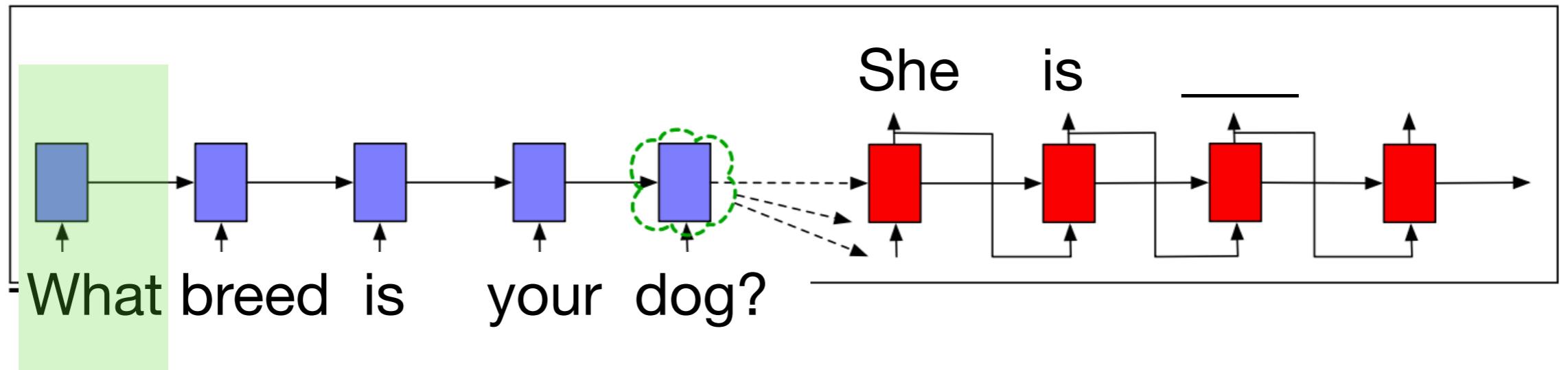
To complete this sentence, we probably want to pay **attention** to certain words in the input.

# Attention: Which inputs matter?



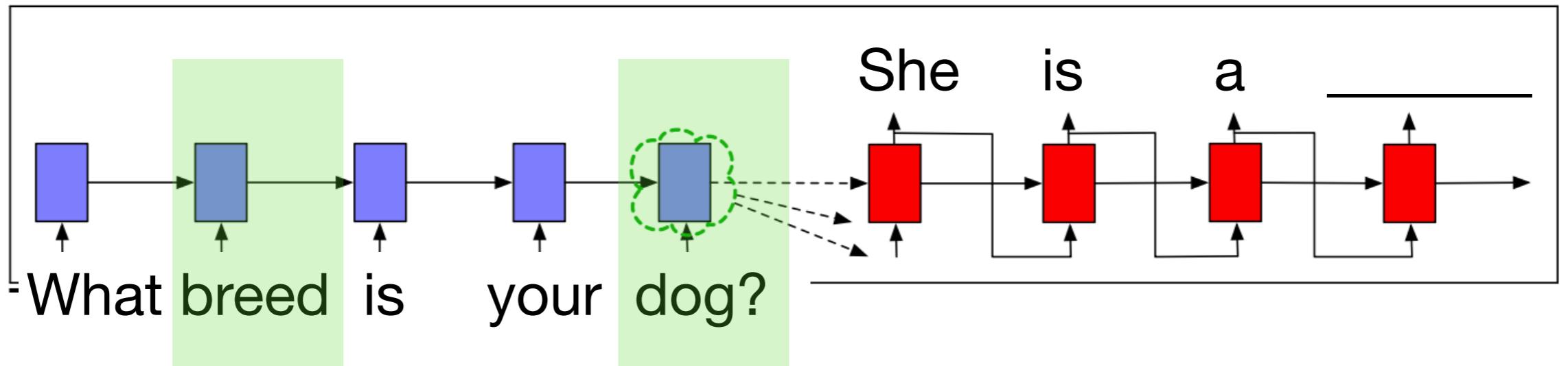
To complete this sentence, we probably want to pay **attention** to certain words in the input.

# Attention: Which inputs matter?



To complete this sentence, we probably want to pay **attention** to certain words in the input.

# Attention: Which inputs matter?



To complete this sentence, we probably want to pay **attention** to certain words in the input.

# From transformers to GPT

## Improving Language Understanding by Generative Pre-Training

Alec Radford  
OpenAI  
alec@openai.com

Karthik Narasimhan  
OpenAI  
karthikn@openai.com

Tim Salimans  
OpenAI  
tim@openai.com

Ilya Sutskever  
OpenAI  
ilyasu@openai.com

### Abstract

Natural language understanding comprises a wide range of diverse tasks such as textual entailment, question answering, semantic similarity assessment, and document classification. Although large unlabeled text corpora are abundant, labeled data for learning these specific tasks is scarce, making it challenging for discriminatively trained models to perform adequately. We demonstrate that large gains on these tasks can be realized by *generative pre-training* of a language model on a diverse corpus of unlabeled text, followed by *discriminative fine-tuning* on each specific task. In contrast to previous approaches, we make use of task-aware input transformations during fine-tuning to achieve effective transfer while requiring minimal changes to the model architecture. We demonstrate the effectiveness of our approach on a wide range of benchmarks for natural language understanding. Our general task-agnostic model outperforms discriminatively trained models that use architectures specifically crafted for each task, significantly improving upon the state of the art in 9 out of the 12 tasks studied. For instance, we achieve absolute improvements of 8.9% on commonsense reasoning (Stories Cloze Test), 5.7% on question answering (RACE), and 1.5% on textual entailment (MultiNLI).

**Key idea:**  
**Pairing transformers with  
unsupervised pre-training**

=> **Unprecedented  
performance**

# From transformers to GPT

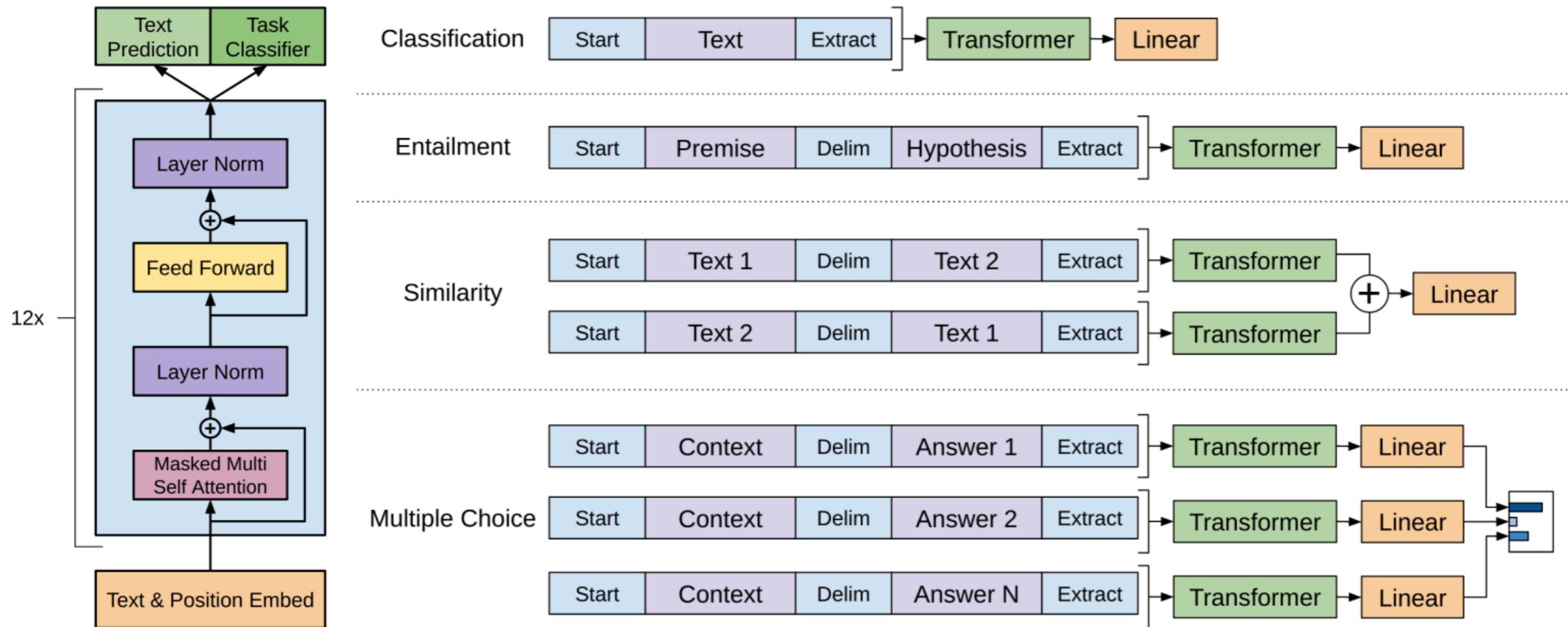


Figure 1: **(left)** Transformer architecture and training objectives used in this work. **(right)** Input transformations for fine-tuning on different tasks. We convert all structured inputs into token sequences to be processed by our pre-trained model, followed by a linear+softmax layer.

Turns out GPT straight out of the box gives funny answers...

How do we make it respond more like a human?

**Prompt:**

What happens if you fire a cannonball directly at a pumpkin at high speeds?

---

**GPT-3 175B completion:**

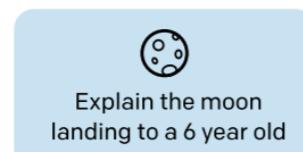
The pumpkin will pull the cannonball in, and the cannonball will bounce off of the pumpkin. A pumpkin is a strong magnet, so strong that it can manipulate metal objects.

# Training language models to follow instructions with human feedback

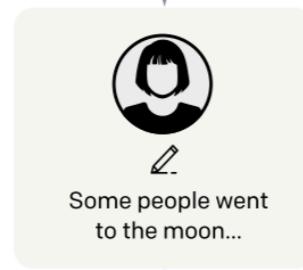
Step 1

**Collect demonstration data, and train a supervised policy.**

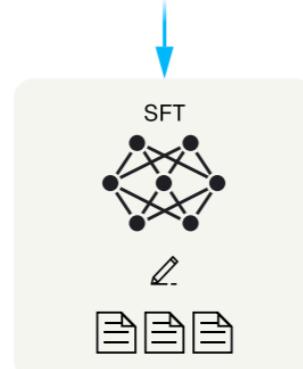
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



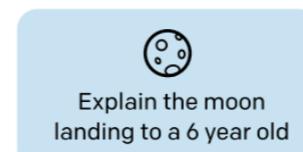
This data is used to fine-tune GPT-3 with supervised learning.



Step 2

**Collect comparison data, and train a reward model.**

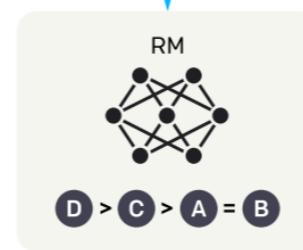
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



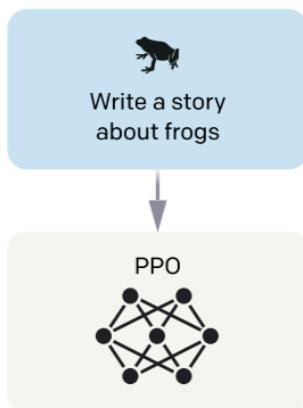
Step 3

**Optimize a policy against the reward model using reinforcement learning.**

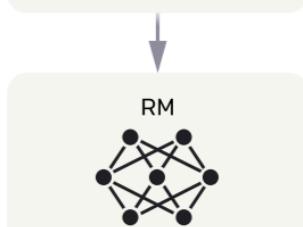
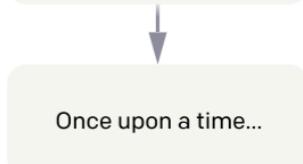
A new prompt is sampled from the dataset.



The policy generates an output.



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.

$r_k$

# Training language models to follow instructions with human feedback

## Prompt:

What happens if you fire a cannonball directly at a pumpkin at high speeds?

---

### GPT-3 175B completion:

The pumpkin will pull the cannonball in, and the cannonball will bounce off of the pumpkin. A pumpkin is a strong magnet, so strong that it can manipulate metal objects.

### InstructGPT 175B completion:

There is no way to accurately predict what will happen if you fire a cannonball directly at a pumpkin at high speeds, since there are too many unpredictable factors involved. Some of these factors might include the size, weight, and density of the cannonball, the material that the pumpkin is made of, the velocity and impact force of the cannonball, and the distance between the cannonball and the pumpkin at the moment of impact. Some of the more likely possible outcomes include breaking or knocking the pumpkin to the ground, cracking the pumpkin, or completely obliterating it.

*Much better! :)*

# Foundation models

## On the Opportunities and Risks of Foundation Models

Rishi Bommasani\* Drew A. Hudson Ehsan Adeli Russ Altman Simran Arora  
Sydney von Arx Michael S. Bernstein Jeannette Bohg Antoine Bosselut Emma Brunskill  
Erik Brynjolfsson Shyamal Buch Dallas Card Rodrigo Castellon Niladri Chatterji  
Annie Chen Kathleen Creel Jared Quincy Davis Dorottya Demszky Chris Donahue  
Moussa Doumbouya Esin Durmus Stefano Ermon John Etchemendy Kawin Ethayarajh  
Li Fei-Fei Chelsea Finn Trevor Gale Lauren Gillespie Karan Goel Noah Goodman  
Shelby Grossman Neel Guha Tatsunori Hashimoto Peter Henderson John Hewitt  
Daniel E. Ho Jenny Hong Kyle Hsu Jing Huang Thomas Icard Saahil Jain  
Dan Jurafsky Pratyusha Kalluri Siddharth Karamcheti Geoff Keeling Fereshte Khani  
Omar Khattab Pang Wei Koh Mark Krass Ranjay Krishna Rohith Kuditipudi  
Ananya Kumar Faisal Ladhak Mina Lee Tony Lee Jure Leskovec Isabelle Levent  
Xiang Lisa Li Xuechen Li Tengyu Ma Ali Malik Christopher D. Manning  
Suvir Mirchandani Eric Mitchell Zanele Munyikwa Suraj Nair Avanika Narayan  
Deepak Narayanan Ben Newman Allen Nie Juan Carlos Niebles Hamed Nilforoshan  
Julian Nyarko Giray Ogut Laurel Orr Isabel Papadimitriou Joon Sung Park Chris Piech  
Eva Portelance Christopher Potts Aditi Raghunathan Rob Reich Hongyu Ren  
Frieda Rong Yusuf Roohani Camilo Ruiz Jack Ryan Christopher Ré Dorsa Sadigh  
Shiori Sagawa Keshav Santhanam Andy Shih Krishnan Srinivasan Alex Tamkin  
Rohan Taori Armin W. Thomas Florian Tramèr Rose E. Wang William Wang Bohan Wu  
Jiajun Wu Yuhuai Wu Sang Michael Xie Michihiro Yasunaga Jiaxuan You Matei Zaharia  
Michael Zhang Tianyi Zhang Xikun Zhang Yuhui Zhang Lucia Zheng Kaitlyn Zhou  
Percy Liang<sup>\*1</sup>

Center for Research on Foundation Models (CRFM)  
Stanford Institute for Human-Centered Artificial Intelligence (HAI)  
Stanford University

This paper coined the term  
“Foundation Models”

*AI is undergoing a paradigm shift with the rise of models (e.g., BERT, DALL-E, GPT-3) trained on broad data (generally using self-supervision at scale) that can be adapted to a wide range of downstream tasks. We call these models foundation models to underscore their critically central yet incomplete character. This report provides a thorough account of the opportunities and risks of foundation models, ranging from their capabilities (e.g., language, vision, robotic manipulation, reasoning, human interaction) and technical principles (e.g., model architectures, training procedures, data, systems, security, evaluation, theory) to their applications (e.g., law, healthcare, education) and societal impact (e.g., inequity, misuse, economic and environmental impact, legal and ethical considerations). Though foundation models are based on standard deep learning and transfer learning, their scale results in new emergent capabilities, and their effectiveness across so many tasks incentivizes homogenization. Homogenization provides powerful leverage but demands caution, as the defects of the foundation model are inherited by all the adapted models downstream. Despite the impending widespread deployment of foundation models, we currently lack a clear understanding of how they work, when they fail, and what they are even capable of due to their emergent properties. To tackle these questions, we believe much of the critical research on foundation models will require deep interdisciplinary collaboration commensurate with their fundamentally sociotechnical nature.*

# Foundation models

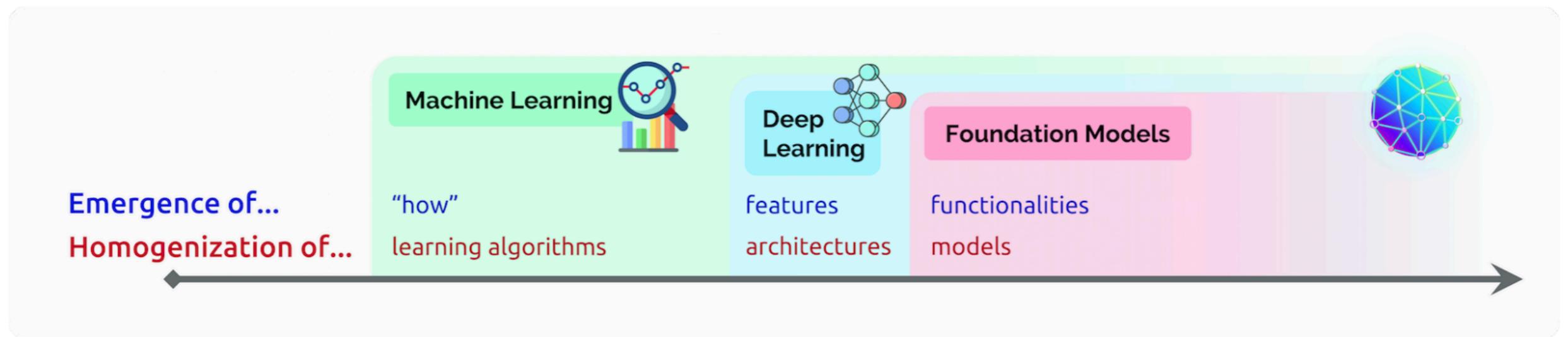
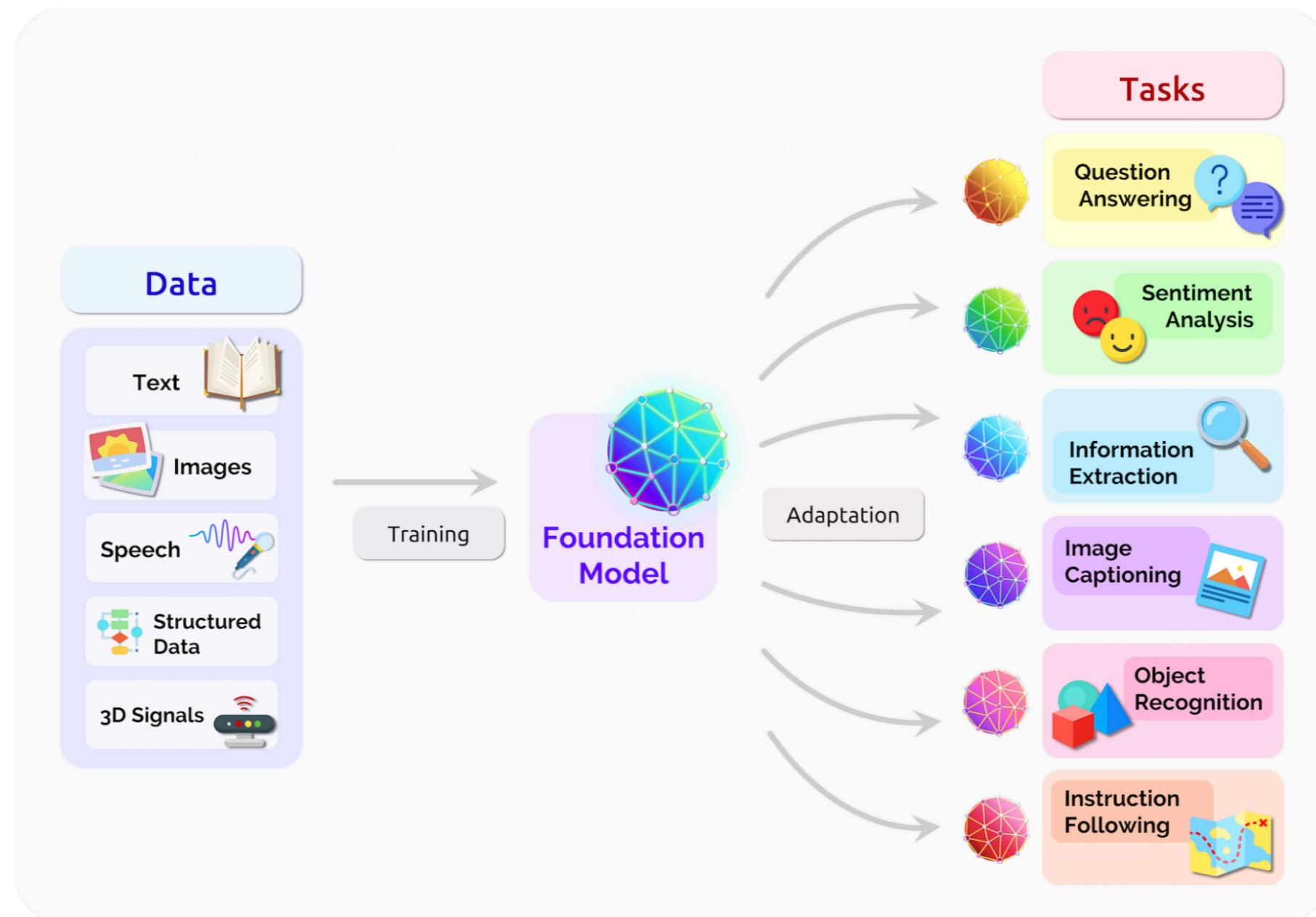


Fig. 1. The story of AI has been one of increasing *emergence* and *homogenization*. With the introduction of machine learning, *how* a task is performed emerges (is inferred automatically) from examples; with deep learning, the high-level features used for prediction emerge; and with foundation models, even advanced functionalities such as in-context learning emerge. At the same time, machine learning homogenizes learning algorithms (e.g., logistic regression), deep learning homogenizes model architectures (e.g., Convolutional Neural Networks), and foundation models homogenizes the model itself (e.g., GPT-3).

# Foundation models



# References

- “Dive into Deep Learning” by Aston Zhang, Alexander J. Smola, Zachary Lipton, Mu Li
- “Speech and Language Processing” by Dan Jurafsky and James H. Martin (<https://web.stanford.edu/~jurafsky/slp3/>)
- [https://cs229.stanford.edu/main\\_notes.pdf](https://cs229.stanford.edu/main_notes.pdf)