



From CS188, FA18, UC Berkeley: <http://ai.berkeley.edu>.



Natural Language Processing

Introduction to Data Science
Spring 1403

Yadollah Yaghoobzadeh

Agenda

- What is NLP?
- Why NLP is important?
- NLP applications (translation, sentiment analysis, summarization)
- Word embedding (word2vec)
- RNN (sequence processing)
- Attention mechanism, Transformers
- Language modeling (task definition)
- From LMs to general-purpose chatbots

Goal

Comprehension and generation of **natural language**

3

Natural language

- Languages that evolved naturally through human use
 - e.g., Spanish, English, Arabic, Hindi, etc.



4

Machine translation

Google Translate

The screenshot shows the Google Translate web interface. At the top, there are tabs for 'Text' and 'Documents'. Below that, language selection bars show 'ENGLISH - DETECTED' on the left and 'PERSIAN' on the right, with other options like GERMAN, FRENCH, PERSIAN, ENGLISH, and SPANISH available. The main area contains a text input field with the English sentence 'We are starting to learn artificial intelligence.' and its Persian translation 'ما در حال پادگیری هوش مصنوعی هستیم.' Below the input field are microphone and speaker icons, and a character count of '49 / 5000'. On the right side, there are sharing and feedback icons, and a link to 'Send feedback'. The entire interface is framed by a light gray border.

Send feedback

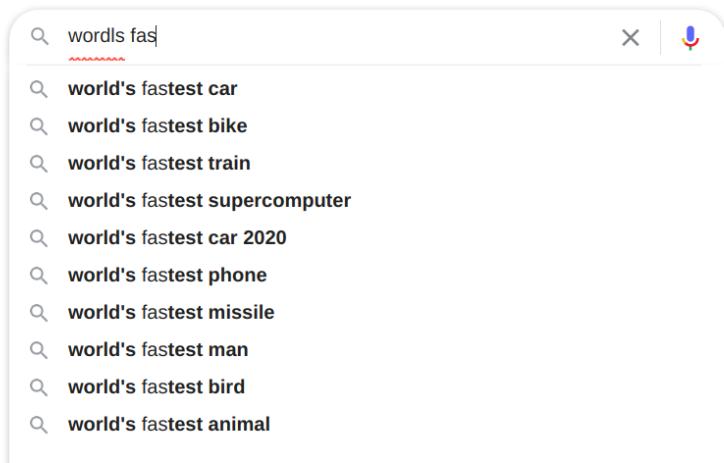
5

Search & QA

The screenshot shows a Google search results page for the query 'how many hours should i sleep'. The search bar at the top contains the query. Below it, the standard Google navigation bar includes 'All', 'Images', 'News', 'Videos', 'Shopping', 'More', 'Settings', and 'Tools'. A snippet of text from the National Sleep Foundation guidelines states: 'National Sleep Foundation guidelines¹ advise that healthy adults need between 7 and 9 hours of sleep per night. Babies, young children, and teens need even more sleep to enable their growth and development. People over 65 should also get 7 to 8 hours per night.' Below this, a link to 'How Much Sleep Do We Really Need? | Sleep Foundation' is shown. Further down, a 'People also ask' section lists questions such as 'Is it OK to get 5 hours of sleep?', 'Is 6 hours sleep enough?', 'How much sleep do you need by age?', and 'Is it okay to sleep 12 hours a day?'. At the bottom, there is a link to 'Sleep Calculator: How Much Sleep Do You Need? - Healthline' and a brief description about sleep needs. The entire page is framed by a light gray border.

6

Search autocorrect and autocomplete



7

Social media analysis



Chatbots

10:05

Hello

Hi, what can we do for you?

Event Feedback

Thanks for coming along today - we'd love to hear what you think about today's event and the presentations. Let's get started.

So, what do you think? Give me you gut feeling!

9

Hiring and recruitment



10

Voice assistants



11

Grammar checkers

The most common type of marketing channel is the wholesale market. Varies kinds of **produce** are supplied from different areas are assembled at one place and sold thru vegetables s

Replace the word **products**

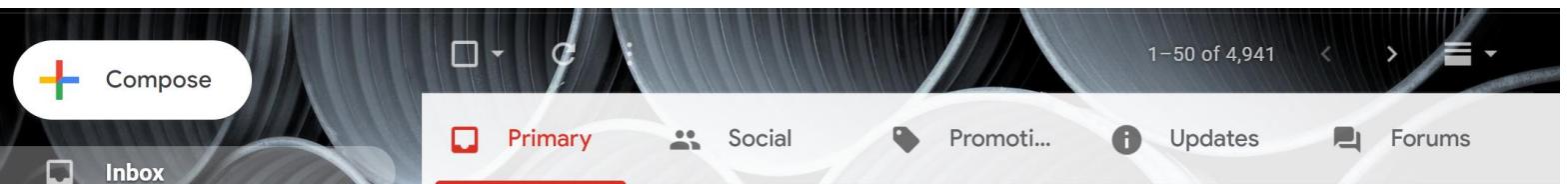
naller regional markets, etc. Fruits and market handling and transport methods.

Dismiss

Suggested by Grammarly

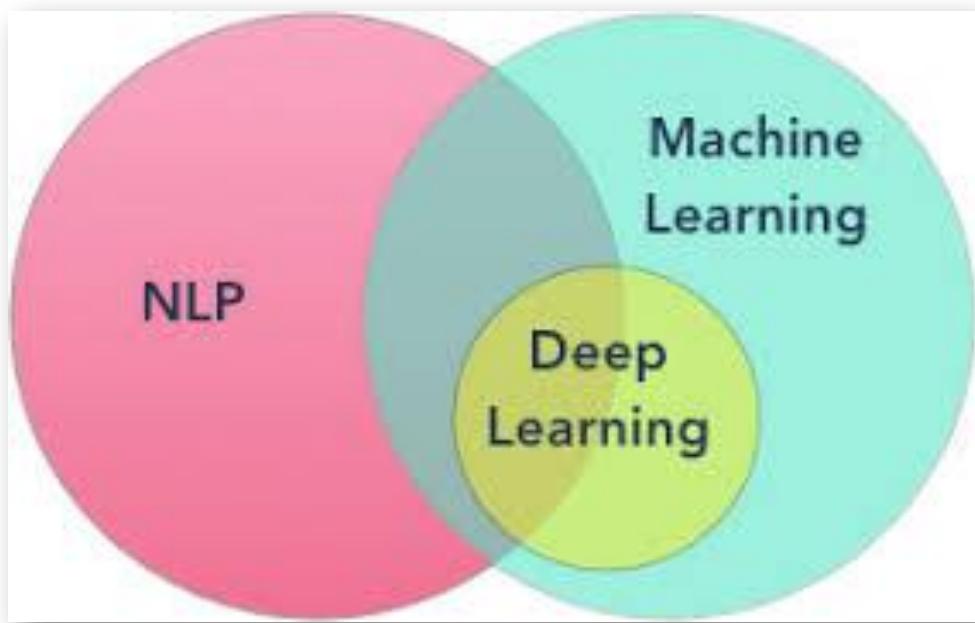
12

Email classification



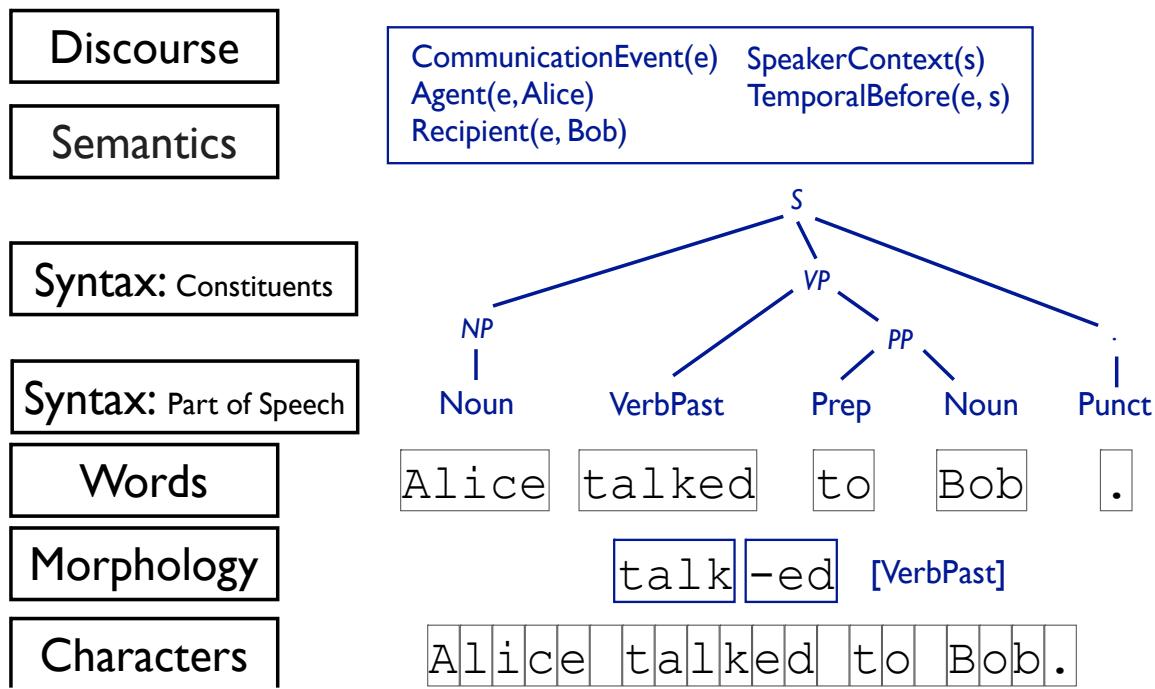
13

NLP is not just machine learning



14

Levels of linguistic structure



15

Deep Learning for Text Classification

Classification

- Output a choice from a fixed set of labels
- For sentiment:
 - Positive/negative
 - Star rating
 - ...

Some examples of binary sentiment classification

this movie was great! would watch again

+

the movie was gross and overwrought, but I liked it

+

this movie was not really very enjoyable

-

Building a classifier

- Let's say we have 10k labeled sentences
- We want to learn a function f that
 - maps an unseen sentence to one of the labels

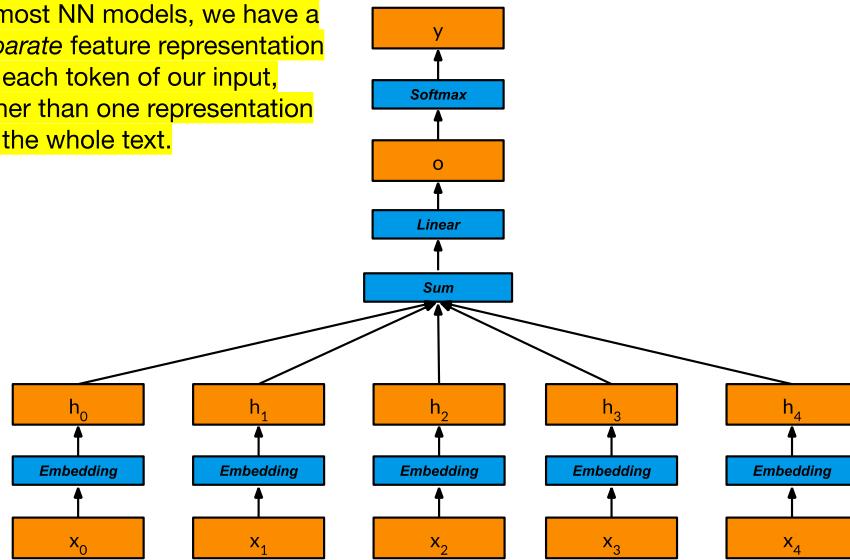
From strings to words

- I don't like any of Ford's trucks.
- I do n't like any of Ford 's trucks .

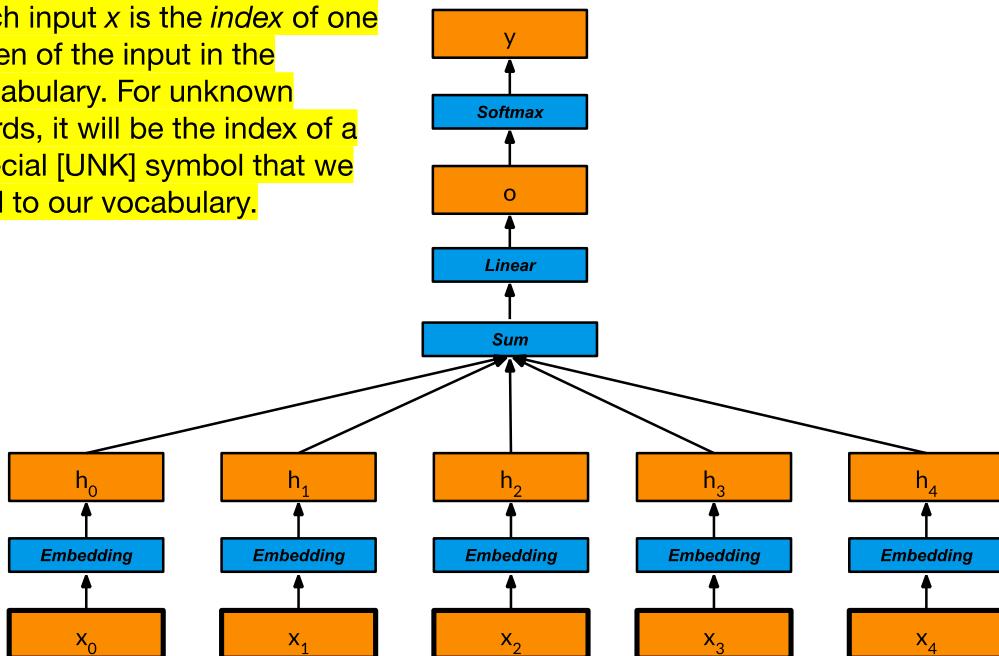
Tokenization: Turning strings into a sequence of symbols (e.g., words, subwords, characters, etc)

NN sentiment classifier

In most NN models, we have a separate feature representation for each token of our input, rather than one representation for the whole text.

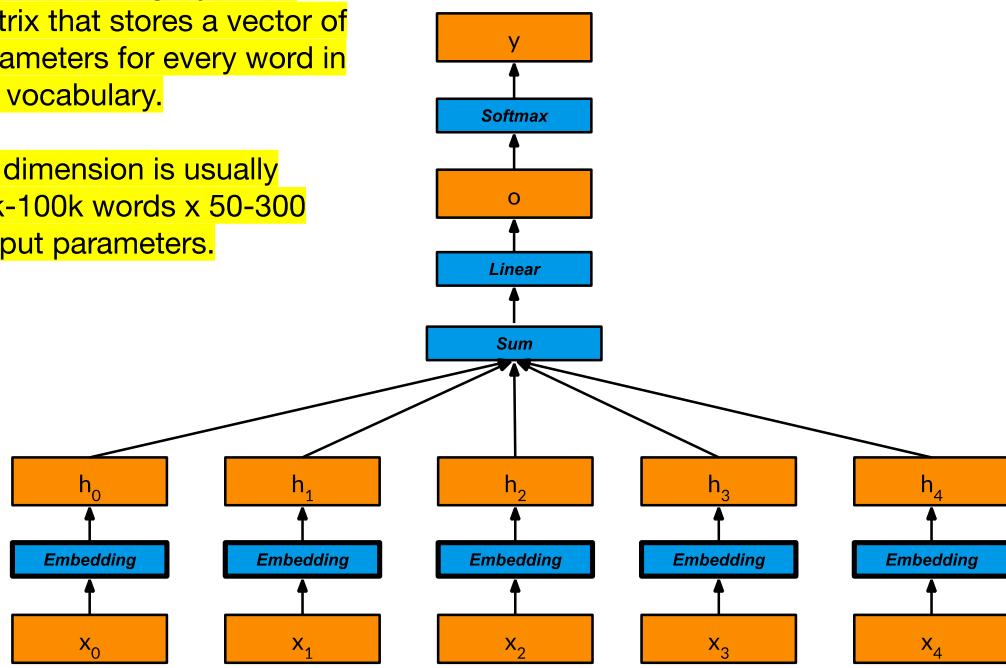


Each input x is the *index* of one token of the input in the vocabulary. For unknown words, it will be the index of a special [UNK] symbol that we add to our vocabulary.

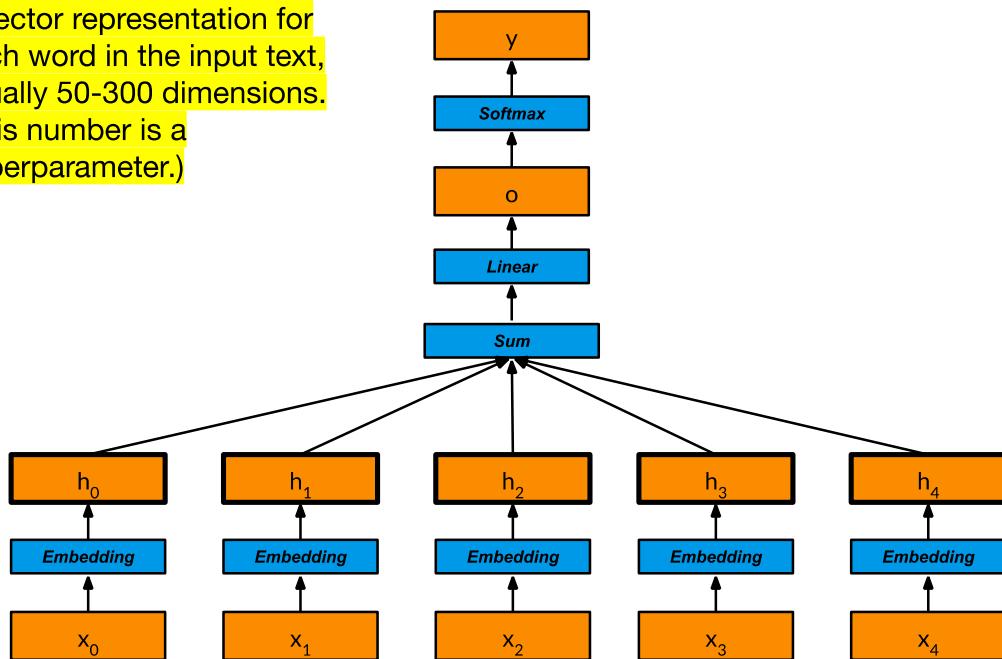


The *embedding layer* is a matrix that stores a vector of parameters for every word in the vocabulary.

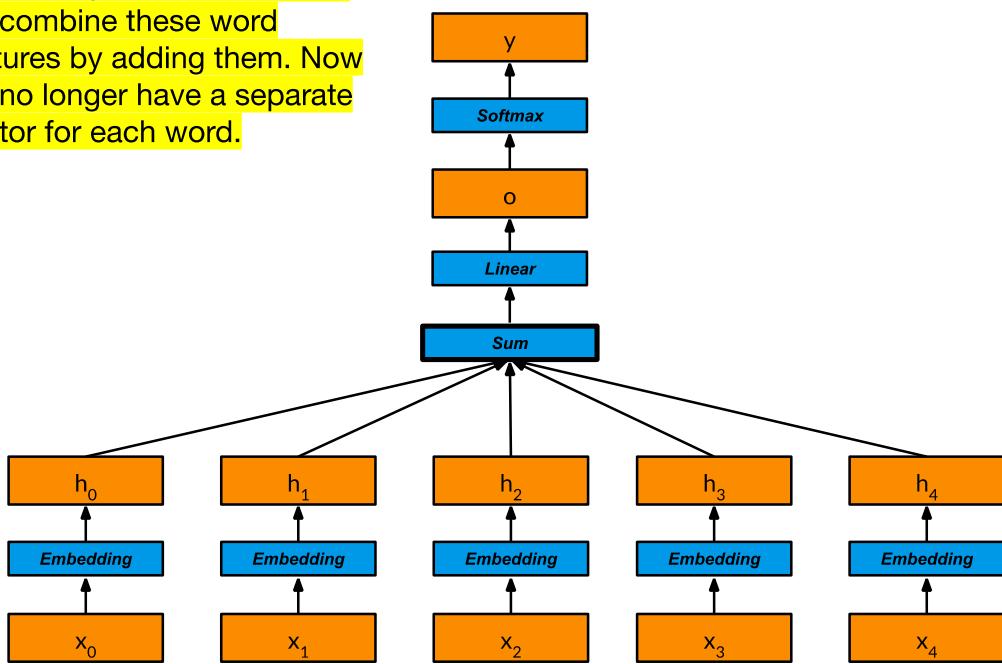
It's dimension is usually 10k-100k words x 50-300 output parameters.



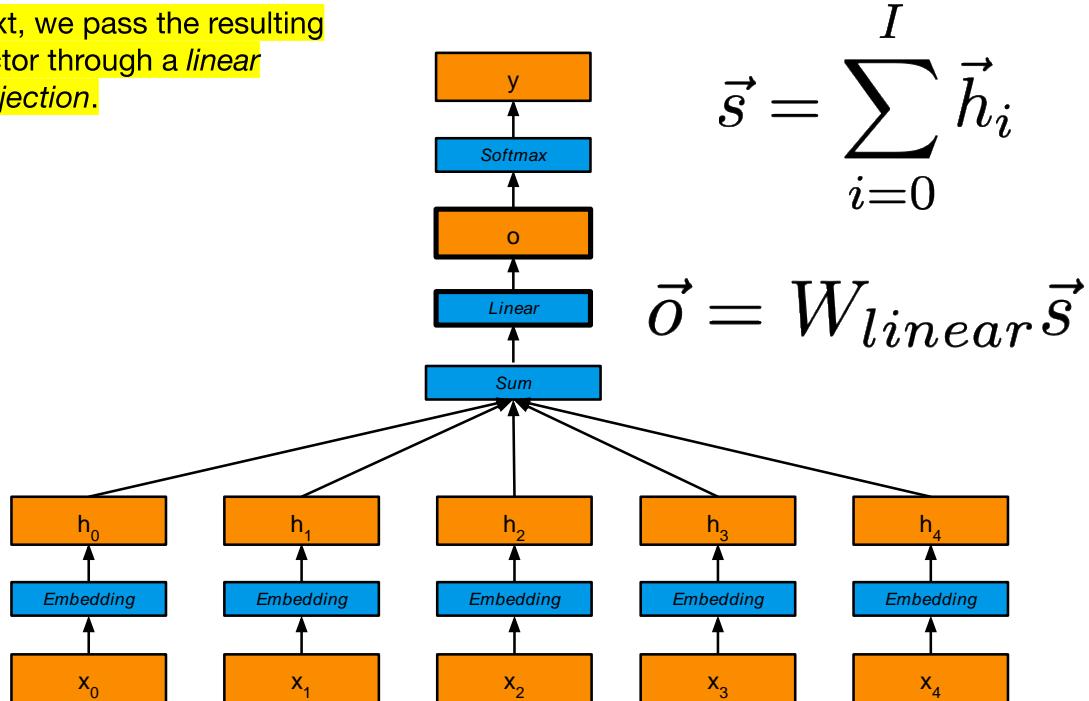
The embedding layer produces a vector representation for each word in the input text, usually 50-300 dimensions. (This number is a hyperparameter.)



In this simple neural network, we combine these word features by adding them. Now we no longer have a separate vector for each word.

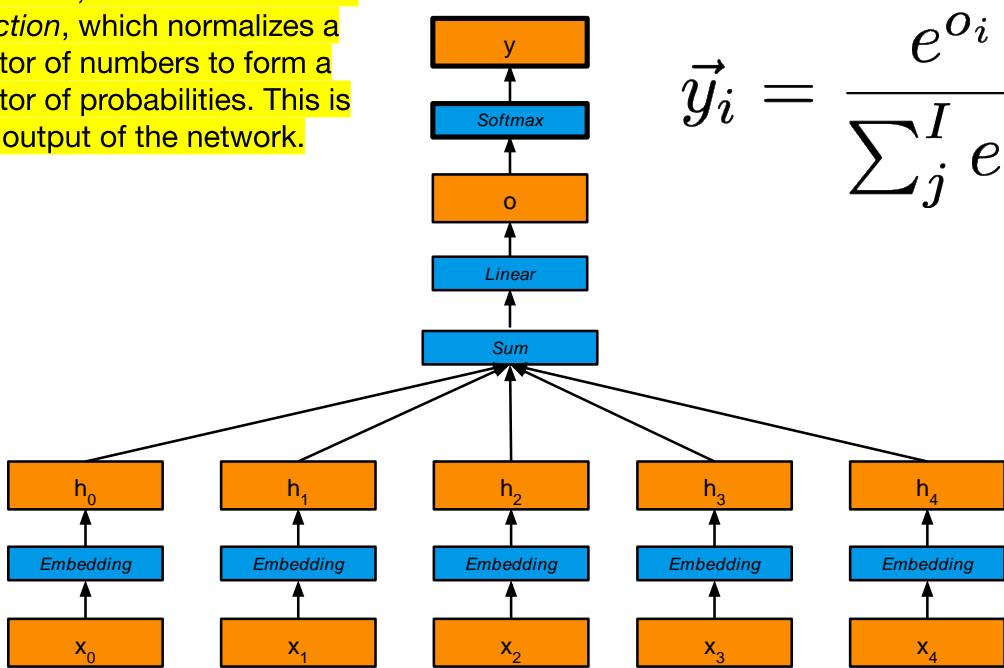


Next, we pass the resulting vector through a *linear projection*.



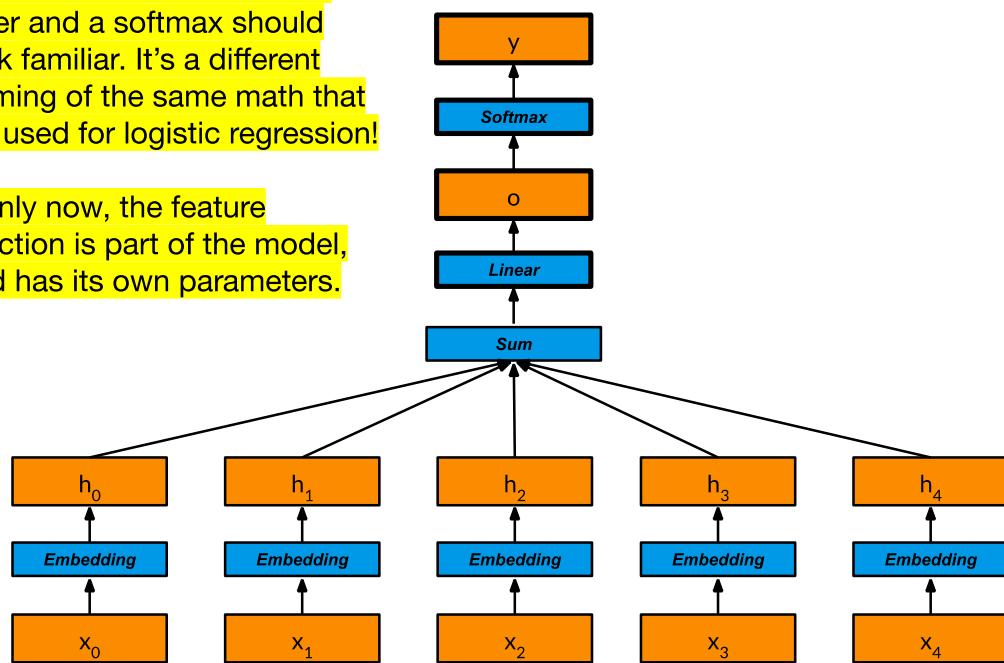
After that, we use the *softmax* function, which normalizes a vector of numbers to form a vector of probabilities. This is the output of the network.

$$\vec{y}_i = \frac{e^{o_i}}{\sum_j e^{o_j}}$$

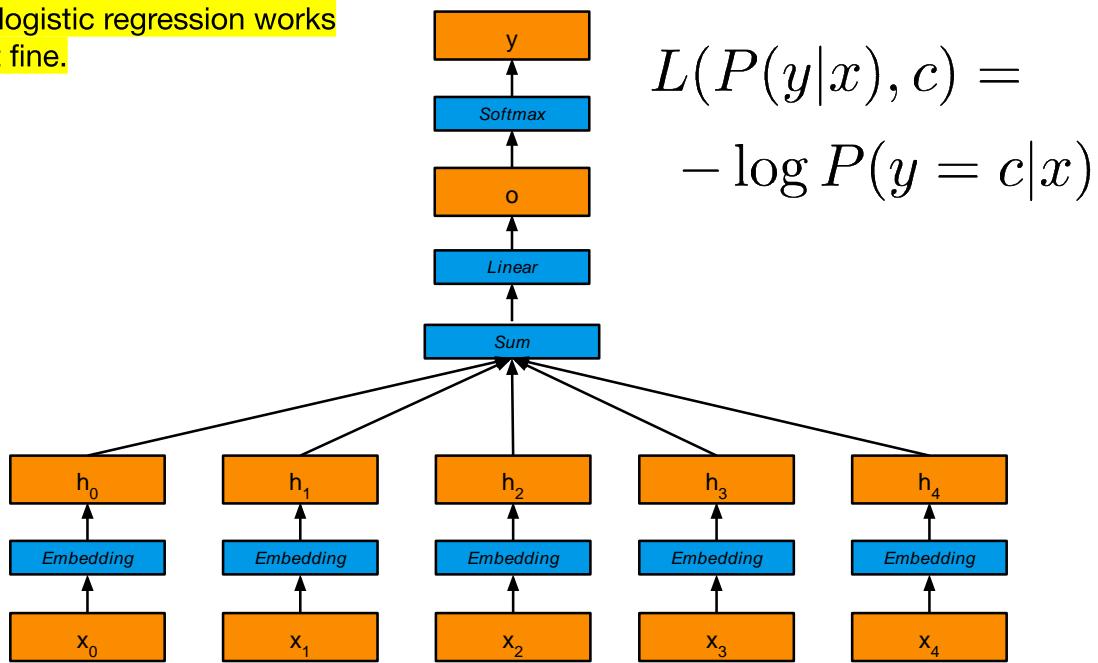


This combination of a linear layer and a softmax should look familiar. It's a different framing of the same math that we used for logistic regression!

...only now, the feature function is part of the model, and has its own parameters.



The *loss function* that we used
for logistic regression works
just fine.

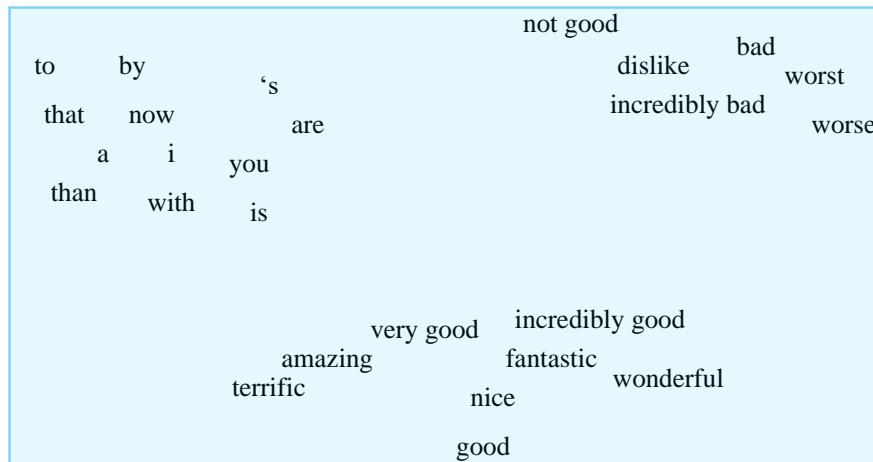


$$L(P(y|x), c) = -\log P(y = c|x)$$

Word embedding

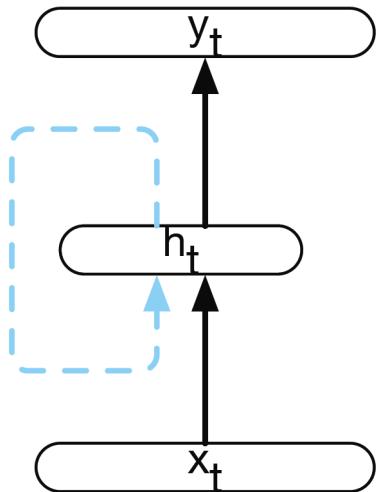
Word embedding

- Each word = a vector
- Similar words are "nearby in space"



RNNs, attention, transformers

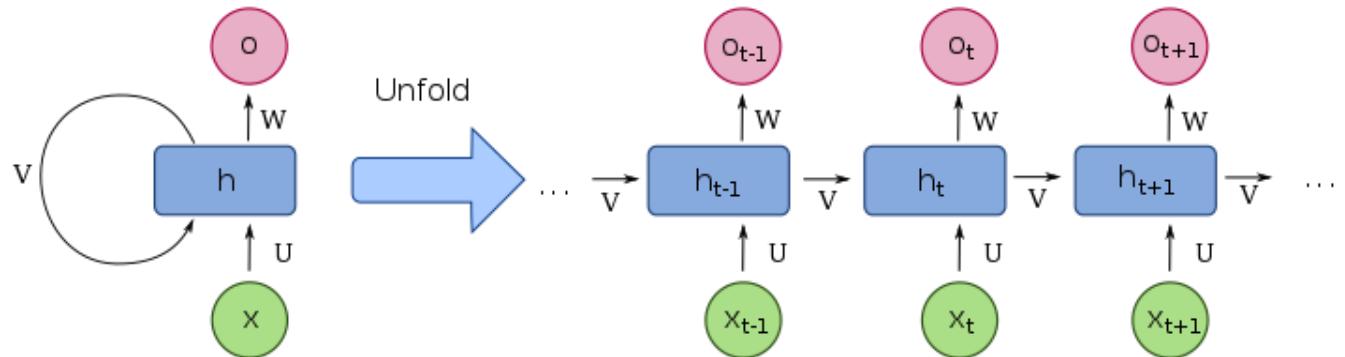
Recurrent Neural Networks (RNNs)



- RNNs take the previous output or hidden states as inputs!
The composite input at time t has some historical information about the happenings at time $T < t$
- RNNs are useful as their intermediate values (state) can store information about past inputs for a time that is not fixed a priori

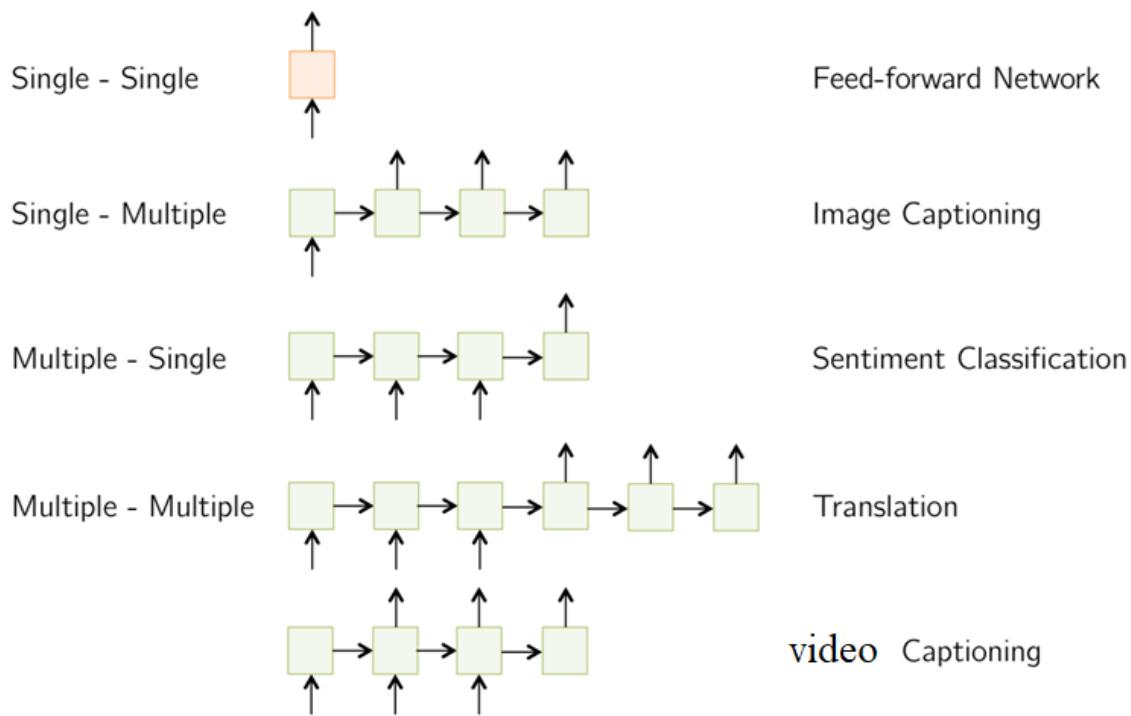
33

Recurrent Neural Networks (RNNs)



34

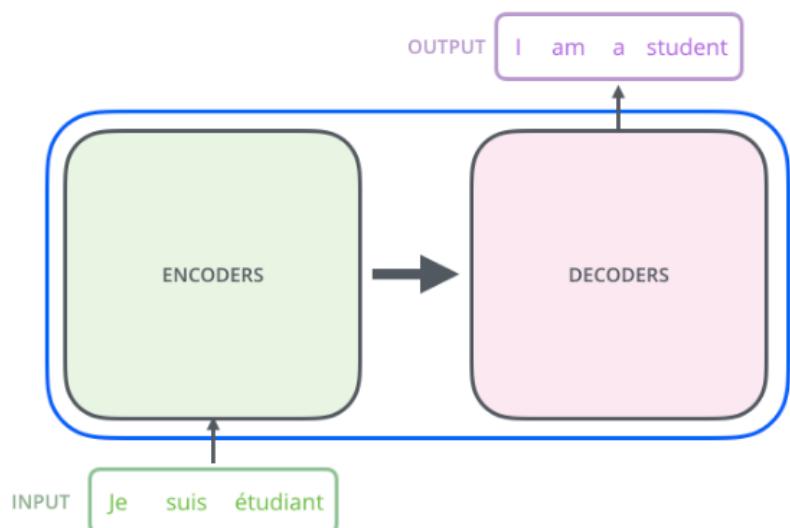
RNNs applications



35

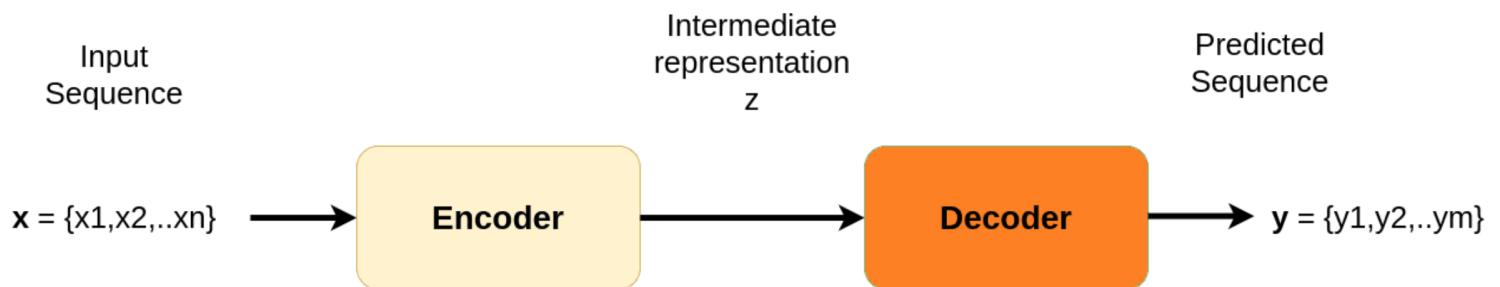
Encoder-decoder networks

- Used in a wide range of applications including machine translation, summarization, question answering, and dialogue modeling.
- RNNs were the most widely-used and successful architecture for both the encoder and decoder.



36

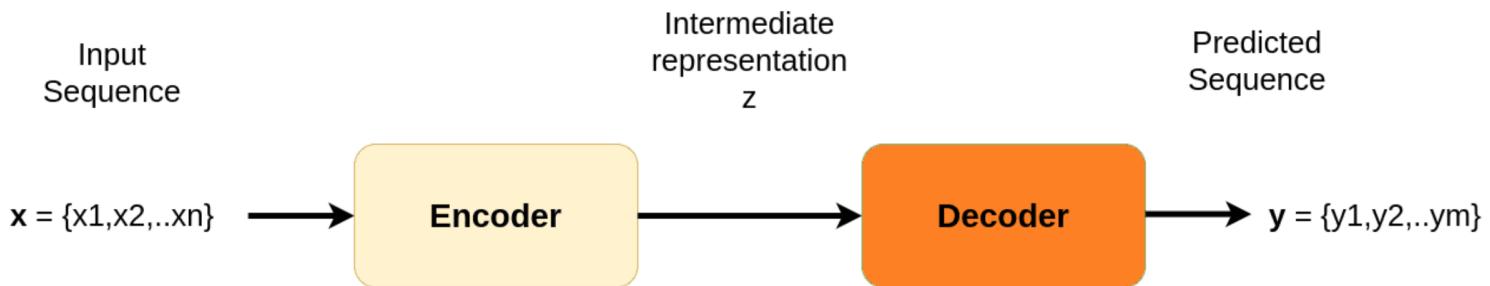
Encoder-decoder: seq2seq



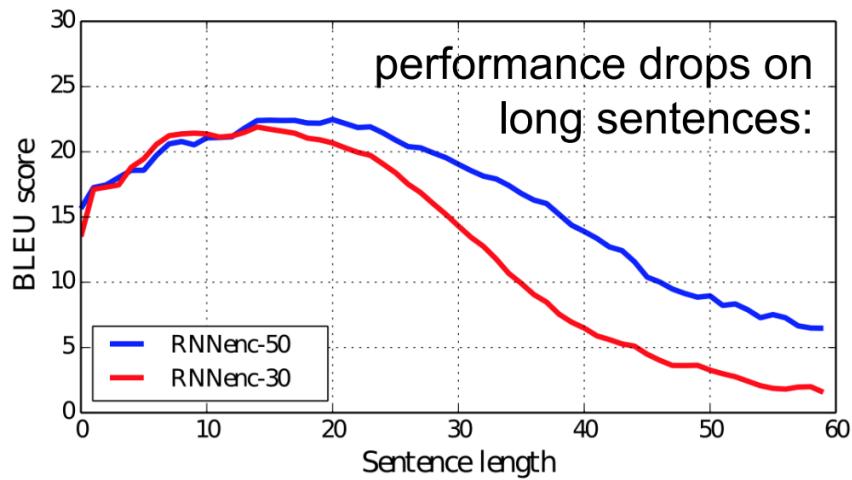
<https://theaisummer.com/attention/>

37

What if sequence length is high (say > 30)?



The vector z needs to capture all the information about the source sentence.

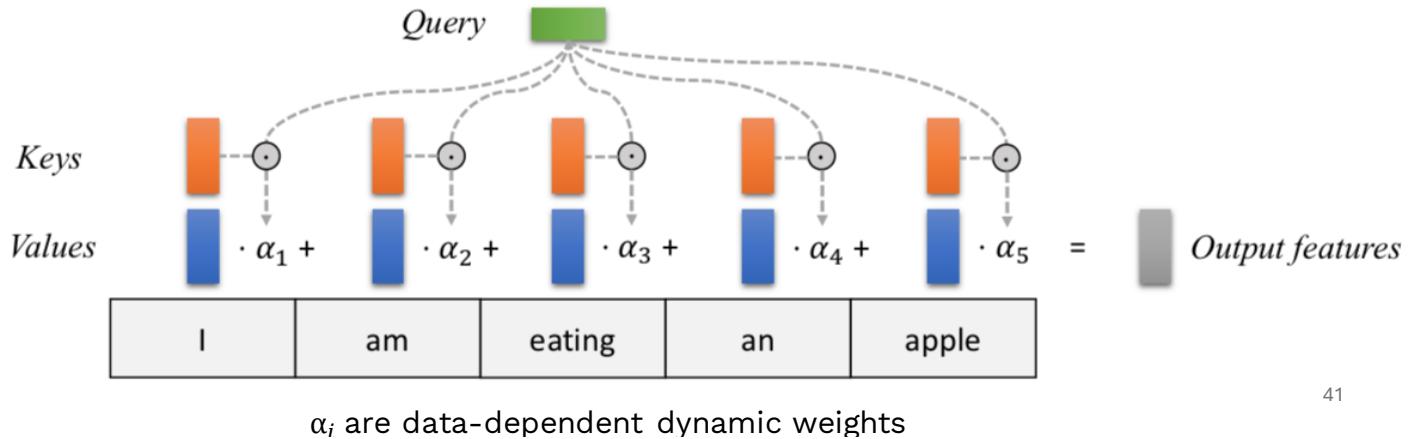


fixed size representation can be the bottleneck

The core idea of **attention** is that the context vector z should have access to **all** parts of the input sequence instead of just the last one.

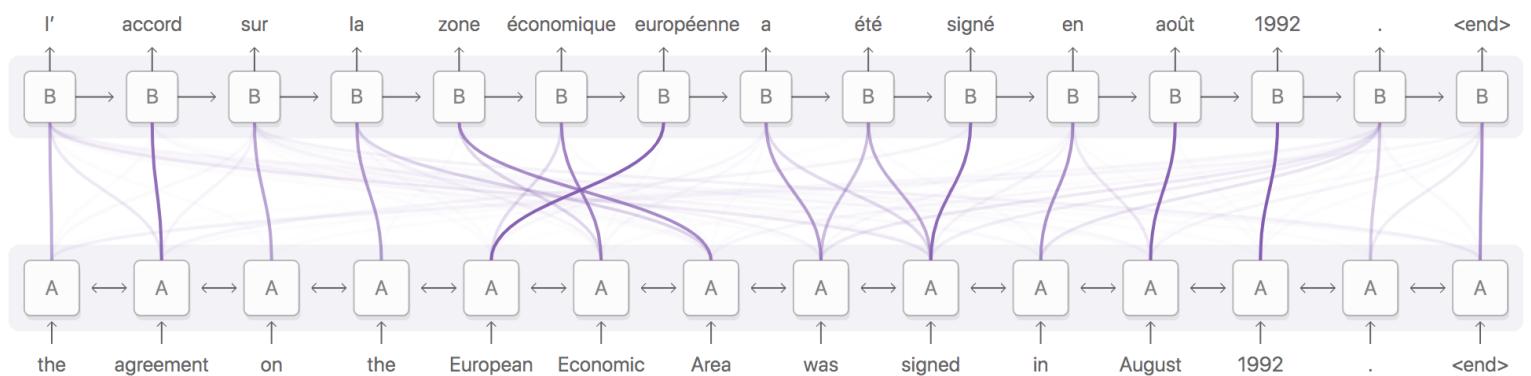
Attention

- A weighted average of (sequence) elements with the weights depending on an input query.
 - The idea and the name taken from the intuition that humans attend to certain things at each time.



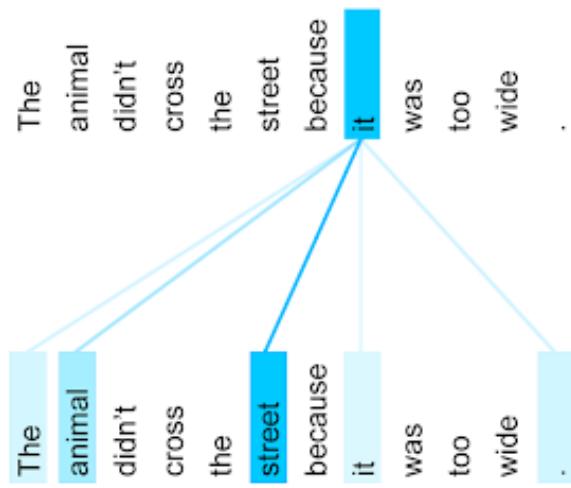
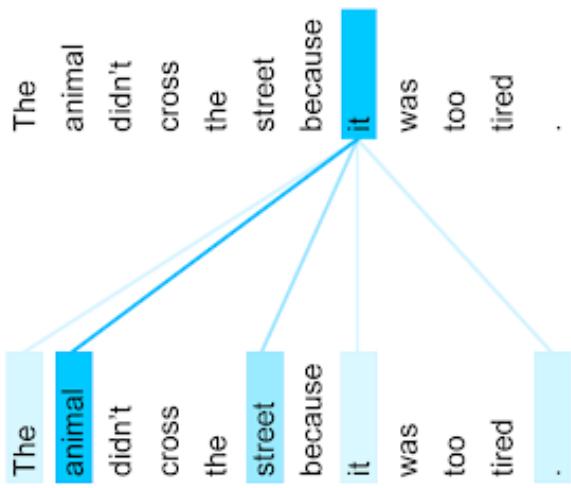
41

MT



Self-Attention

- For each word, self-attention allows the model to look at other positions in the input for a better encoding for this word.

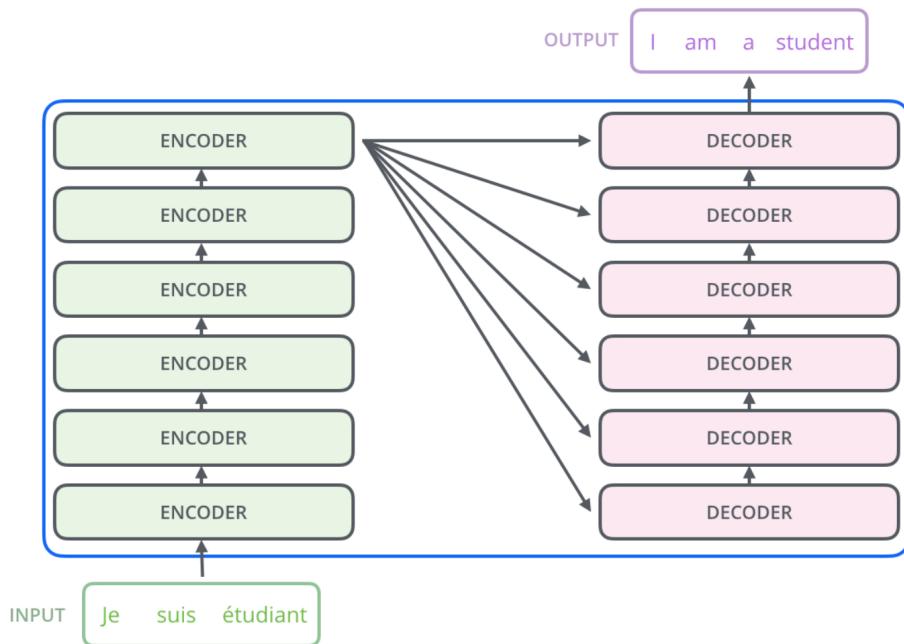


43

Transformer

44

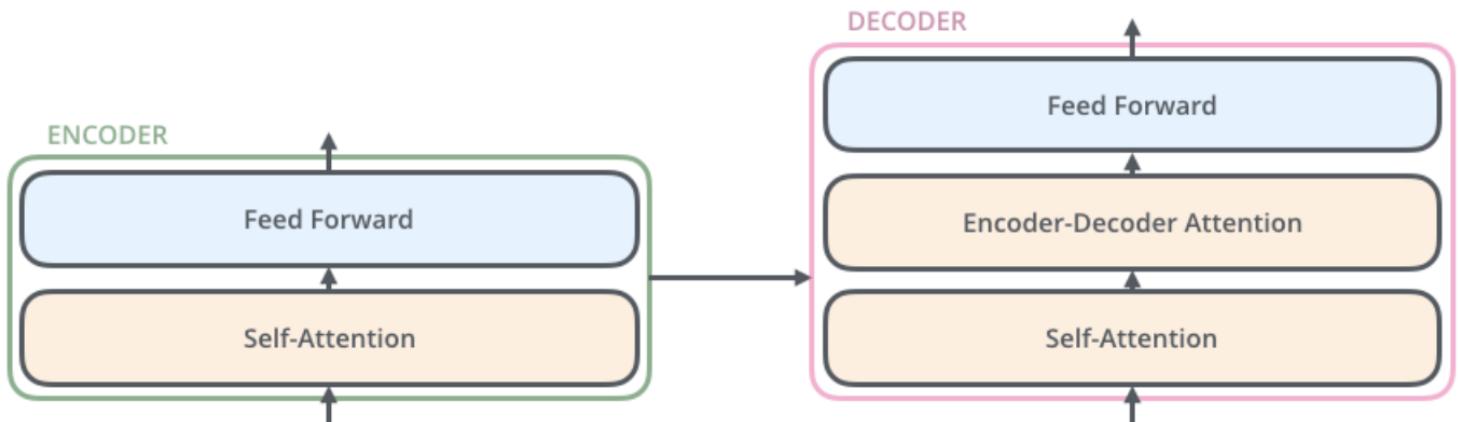
Transformer Model



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

45

Transformer Model

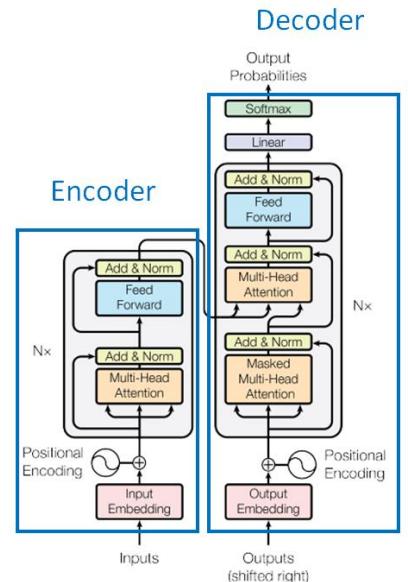


<http://jalammar.github.io/illustrated-transformer/>

46

Transformer model

- Non-recurrent sequence to sequence encoder-decoder model
 - Eliminate recurrence, allows for significantly more parallelization
- Three kinds of attentions:
 - The input and output tokens (solved by traditional attention mechanism)
 - The input tokens themselves
 - The output tokens themselves
- Extend the (self-)attention mechanism to processing input and output sentences as well.



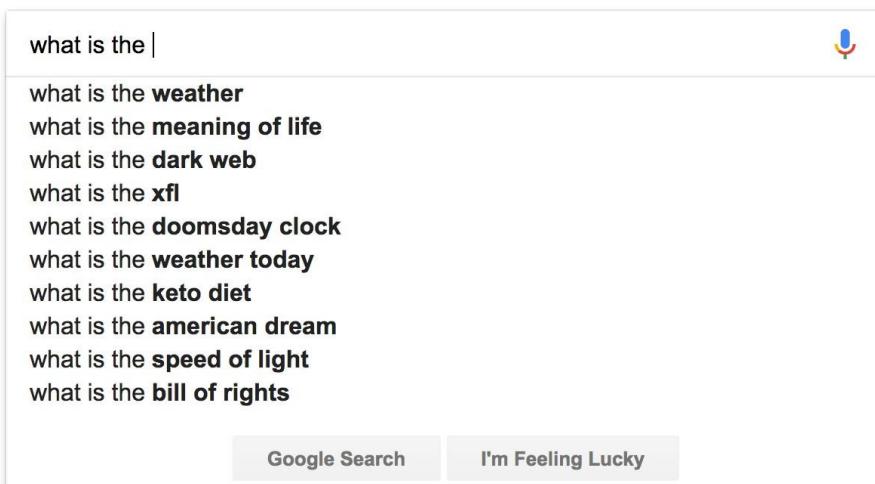
Language modeling

Slides are adopted from Advanced NLP course at UMass.

Intro

- Language models assign a probability to a piece of text.
- Why would we ever want to do this?
 - Translation:
 $P(i \text{ flew to the movies}) <<< P(i \text{ went to the movies})$
 - Speech recognition:
 $P(i \text{ saw a van}) >>> P(\text{eyes awe of an})$

You use Language Models every day!



Probabilistic Language Modeling

- Goal: compute the probability of a sentence or sequence of words:
 - $P(W) = P(w_1, w_2, w_3, w_4, w_5, \dots, w_n)$
- Related task: probability of an upcoming word:
 - $P(w_5|w_1, w_2, w_3, w_4)$
- A model that computes either of these:
 - $P(W)$ or $P(w_n|w_1, w_2, \dots, w_{n-1})$ is called a *language model* or LM

58

How to compute $P(W)$

- How to compute this joint probability:
 - $P(\text{its, water, is, so, transparent, that})$
- Intuition: let's rely on the Chain Rule of Probability

59

Reminder: The Chain Rule

- Recall the definition of conditional probabilities

$$P(B|A) = P(A, B)/P(A) \quad \text{Rewriting: } P(A, B) = P(A)P(B|A)$$

- More variables:

$$P(A, B, C, D) = P(A)P(B|A)P(C|A, B)P(D|A, B, C)$$

- The Chain Rule in General

$$\begin{aligned} & P(x_1, x_2, x_3, \dots, x_n) \\ &= P(x_1)P(x_2|x_1)P(x_3|x_1, x_2) \dots P(x_n|x_1, \dots, x_{n-1}) \end{aligned}$$

60

Chain Rule for computing joint probability

- The Chain Rule applied to compute joint probability of words in sentence

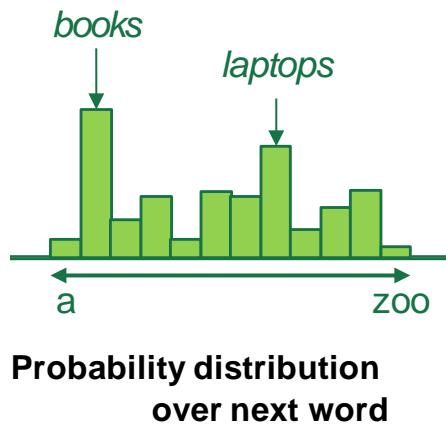
$$P(w_1 w_2 \dots w_n) = \prod_i P(w_i | \underbrace{w_1 w_2 \dots w_{i-1}}_{\text{prefix}})$$

$$\begin{aligned} P(\text{"its water is so transparent"}) &= \\ & P(\text{its}) \times P(\text{water}|\text{its}) \times P(\text{is}|\text{its water}) \\ & \times P(\text{so}|\text{its water is}) \times P(\text{transparent}|\text{its water is so}) \end{aligned}$$

61

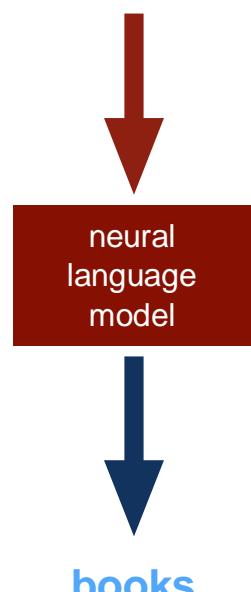
Decoding from an LM

- **Prefix:** “students opened their”



Enter neural networks!

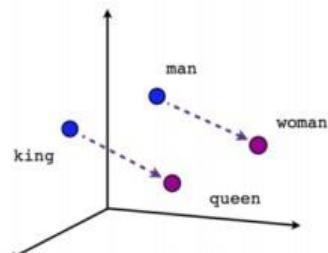
Students opened their



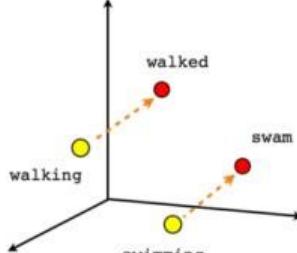
Words as basic building blocks

- Represent words with low-dimensional vectors called embeddings
(Mikolov et al., NIPS 2013)

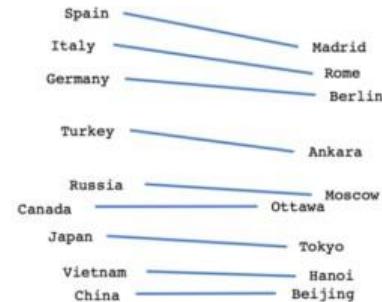
$$\text{king} = [0.23, 1.3, -0.3, 0.43]$$



Male-Female



Verb tense



Country-Capital

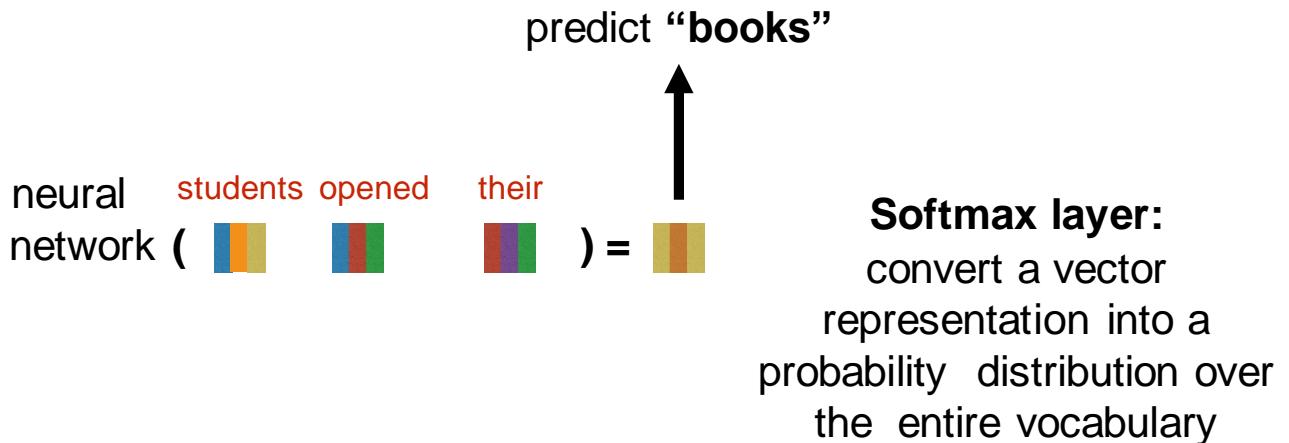
Composing embeddings

- Neural networks **compose** word embeddings into vectors for phrases, sentences, and documents.

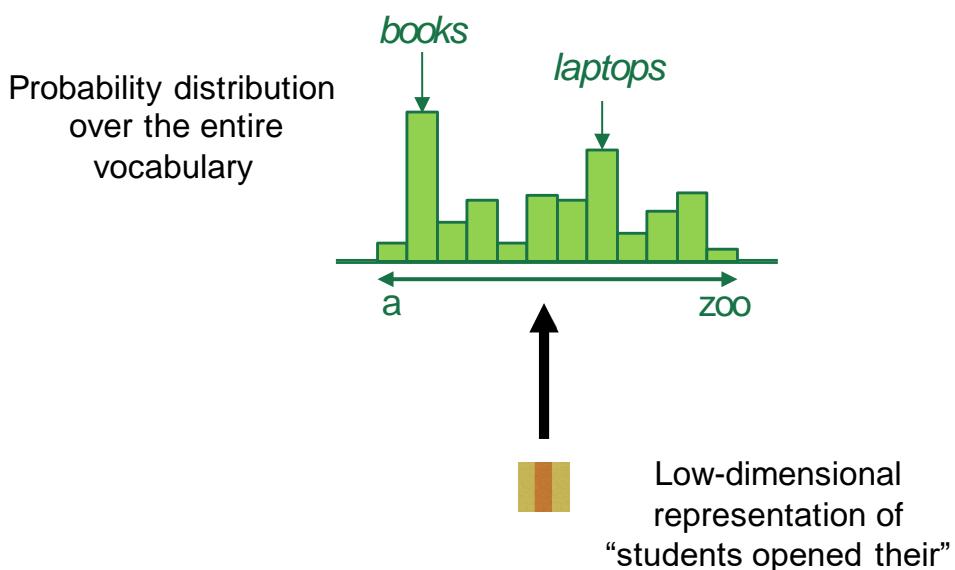
neural students opened their
network (  ) = 

Next word prediction

- Predict the next word from composed prefix representation:



$P(w_i \mid \text{vector for "students opened their"})$



So to sum up...

- Given a d-dimensional vector representation x of a prefix, we do the following to predict the next word:
 - Project it to a V-dimensional vector using a matrix-vector product (a.k.a. a “linear layer”, or a “feedforward layer”), where V is the size of the vocabulary.
 - Apply the softmax function to transform the resulting vector into a probability distribution.

Now that we know how to predict “books”, let’s focus on how to compute the prefix representation x in the first place!



Composition functions

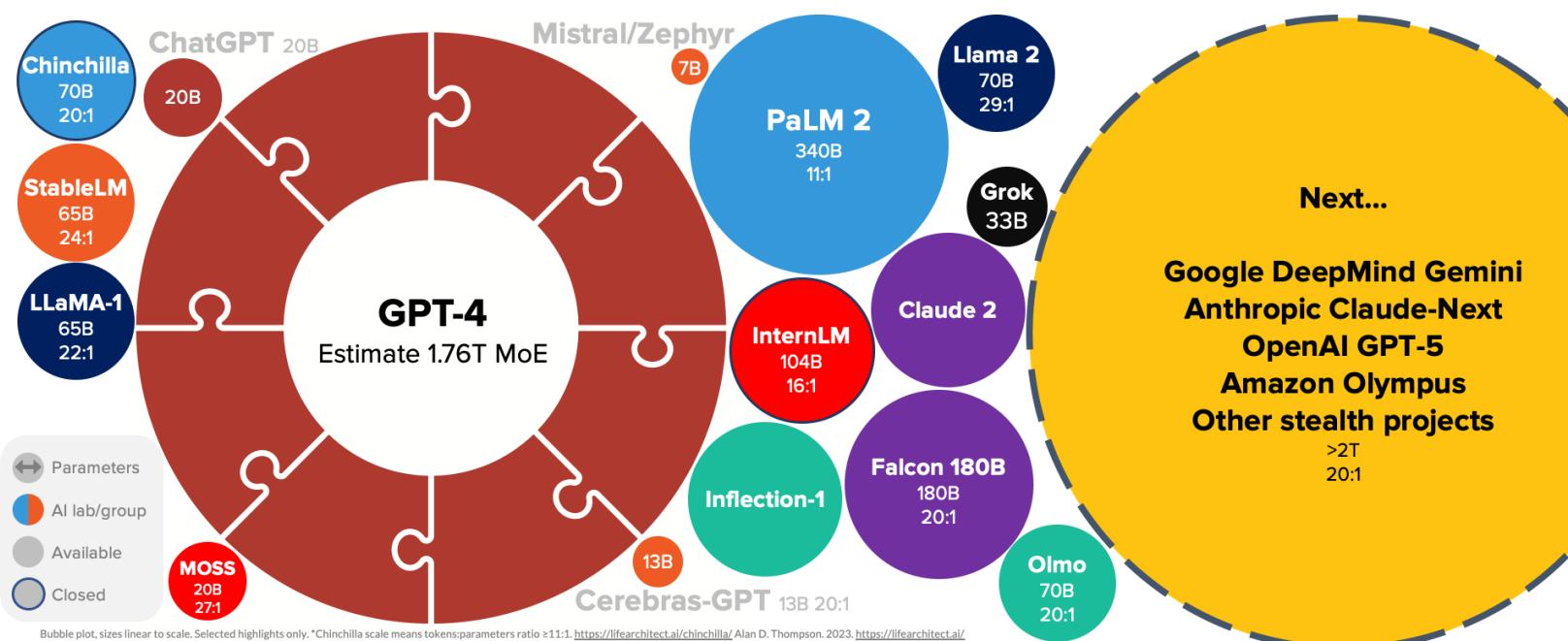
- **Input:** sequence of word embeddings corresponding to the tokens of a given prefix
- **Output:** single vector
- Composition functions:
 - Element-wise functions
 - e.g., just sum up all of the word embeddings!
 - Concatenation
 - Feed-forward neural networks
 - Convolutional neural networks
 - Recurrent neural networks
 - Transformers

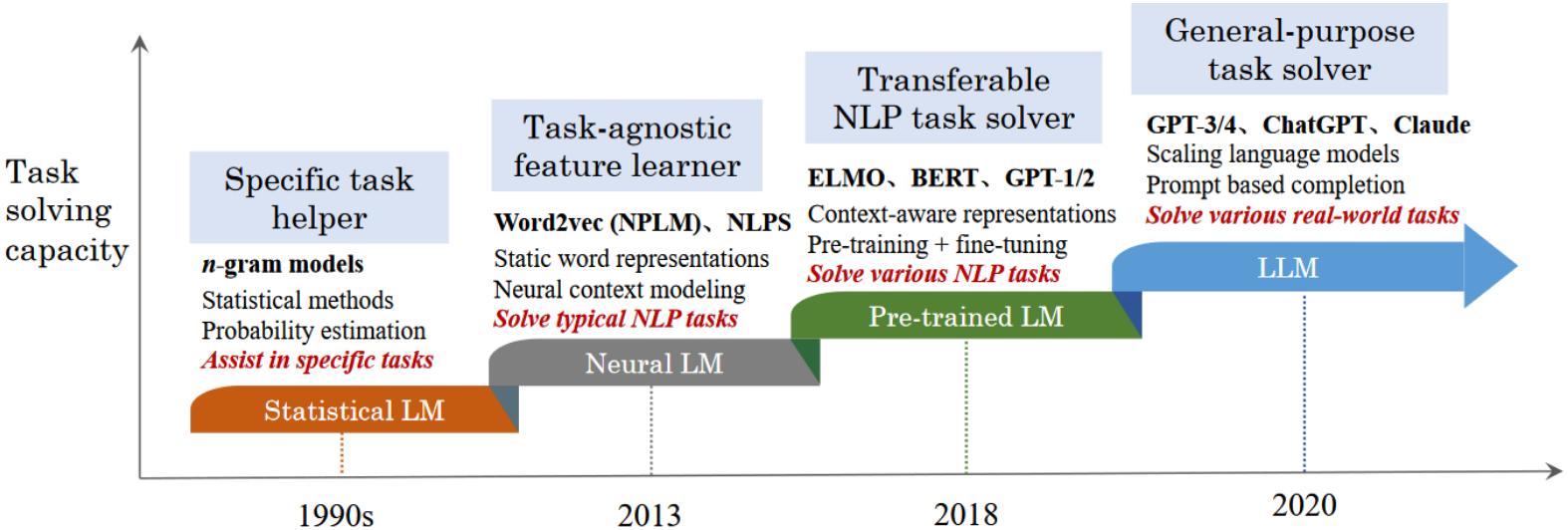
Large language models

What's a large language model?

- ❑ ChatGPT, a chatbot built on top of the GPT LLM released on November 28, 2022, has popularized deep learning models trained with a text corpus.
- ❑ Language models learn a probability distribution over language:
 $P(w_1, \dots, w_m)$
- ❑ Large in terms of training data (e.g., Common Crawl, Wikipedia, GitHub, ...) and parameters (e.g., PaLM has 540 billion parameters; GPT-4 rumored to have 1 trillion).

Language Models (Up to November 2023)

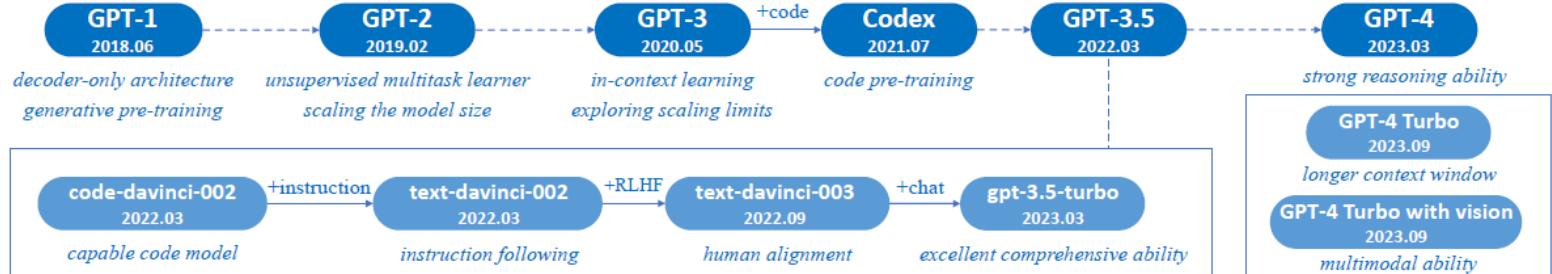




Three kinds of LLM

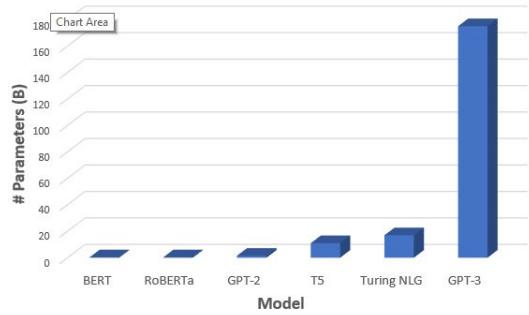
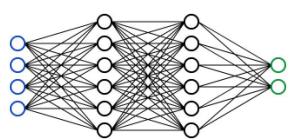
- **Generic language models:** predicting the next token. We will mostly talk about this type today.
- **Instruction tuned**
- **Dialog tuned:** ChatGPT (the base model is hard to interact with).

History of OpenAI GPT models

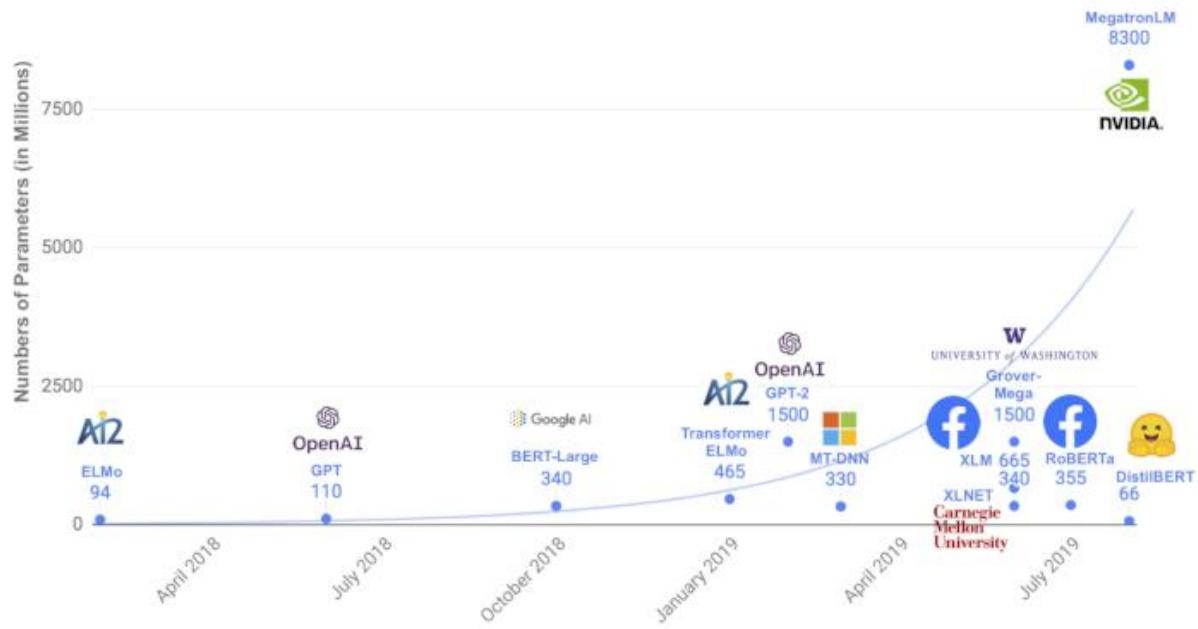


self-supervised learning + scale

In 1885, Stanford University was _____

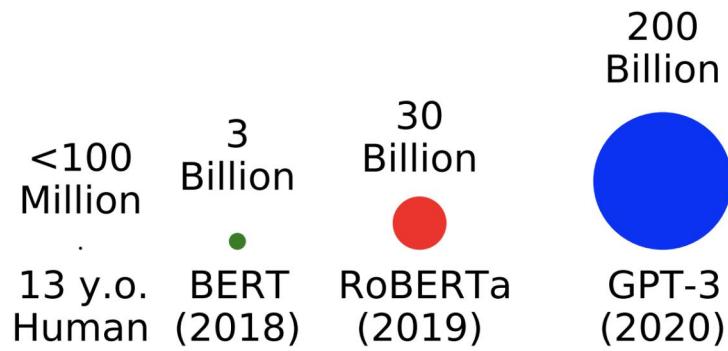


Scale in #parameters



78

Scale in #words



#tokens seen during training

79

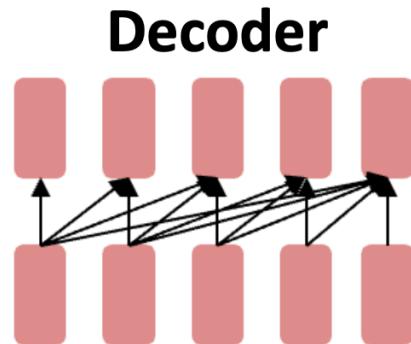
GPTs: Generative Pre-trained Transformers

Large language models

GPT-2 (1.5B parameters; Radford et al., 2019)

Same architecture as GPT, just bigger (117M -> 1.5B)

But trained on much more data: 4GB -> 40GB



 OpenAI

80

Emergent zero-shot learning

The ability to do many tasks with no examples, and no gradient updates

 OpenAI

Language Models are Unsupervised Multitask Learners

Alec Radford *¹ Jeffrey Wu *¹ Rewon Child¹ David Luan¹ Dario Amodei **¹ Ilya Sutskever **¹

81

Emergent abilities of GPT-3 (2020)

- GPT-3 (175B parameters; Brown et al., 2020)
 - Another increase in size (1.5B -> 175B)
 - and data (40GB -> over 600GB)

Language Models are Few-Shot Learners

Tom B. Brown*

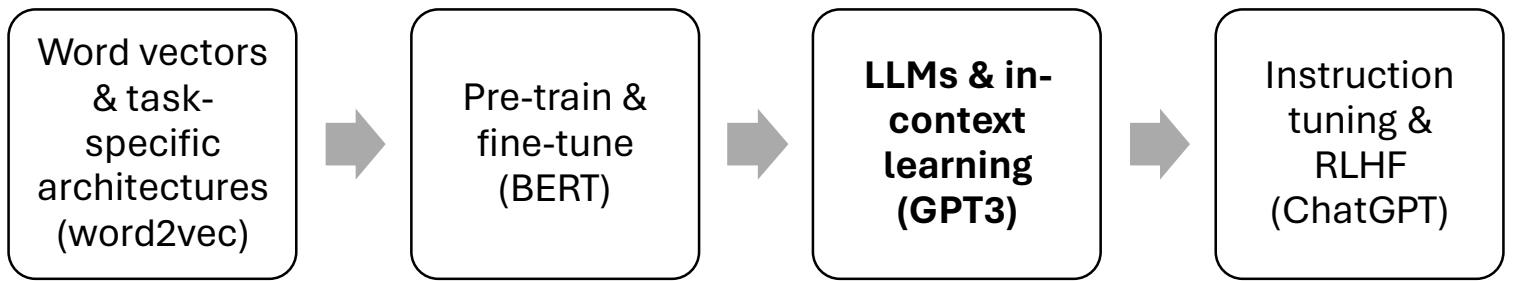
Benjamin Mann*

Nick Ryder*

Melanie Subbiah*

In-context Learning

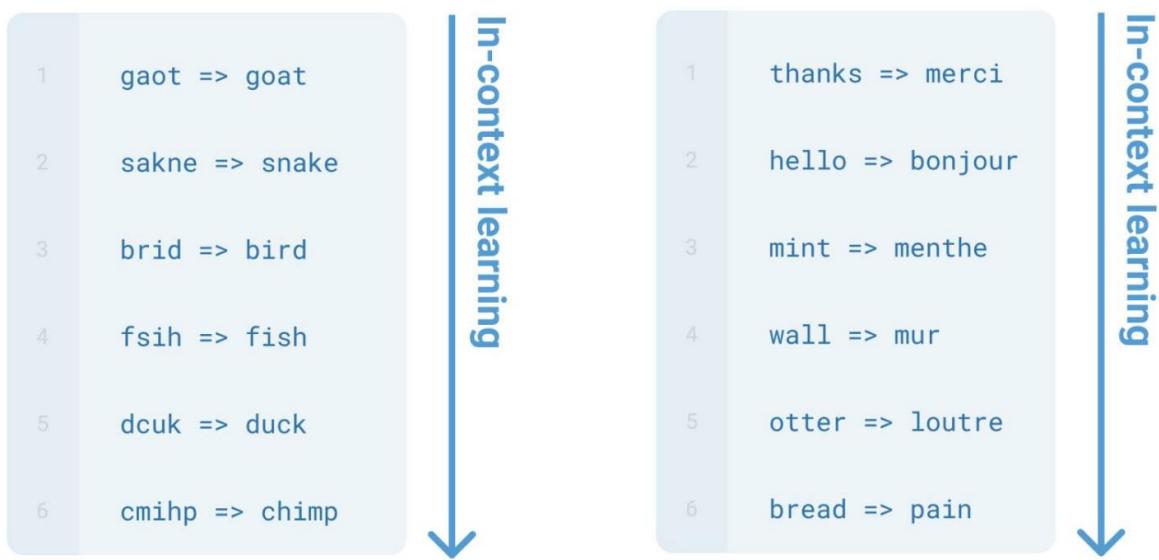
Recent history



84

Emergent in-context learning

 OpenAI



85



Few-shot GPT3 can beat strong task-specific fine-tuned models.

86

Language modeling ≠ assisting users

PROMPT *Explain the moon landing to a 6 year old in a few sentences.*

COMPLETION GPT-3

Explain the theory of gravity to a 6 year old.

Explain the theory of relativity to a 6 year old in a few sentences.

Explain the big bang theory to a 6 year old.

Explain evolution to a 6 year old.

87

Language modeling ≠ assisting users

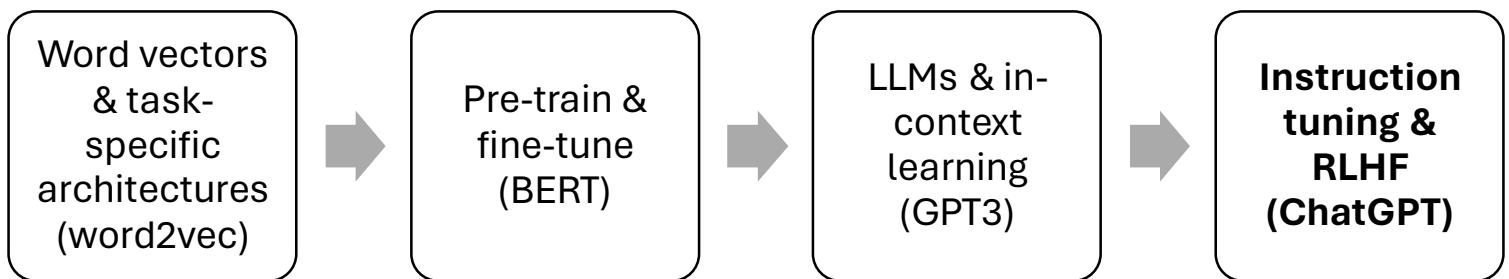
PROMPT *Explain the moon landing to a 6 year old in a few sentences.*

COMPLETION **Human**

A giant rocket ship blasted off from Earth carrying astronauts to the moon. The astronauts landed their spaceship on the moon and walked around exploring the lunar surface. Then they returned safely back to Earth, bringing home moon rocks to show everyone.

Instruction Tuning

Recent history

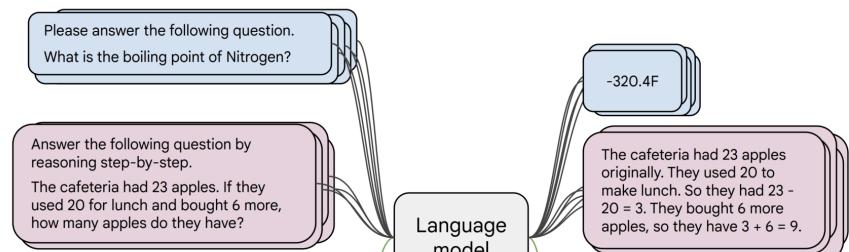


90

Instruction tuning

Collect examples of (instruction, output) pairs across many tasks and finetune an LLM

- Collect examples of (instruction, output) pairs across many tasks and finetune an LM



- Evaluate on unseen tasks

Q: Can Geoffrey Hinton have a conversation with George Washington?
Give the rationale before answering.

Geoffrey Hinton is a British-Canadian computer scientist born in 1947. George Washington died in 1799. Thus, they could not have had a conversation together. So the answer is "no".

[FLAN-T5; Chung et al., 2022]

91

Limitations of instruction-tuning

- No right answer for tasks like open-ended creative generation
- Mismatch between the LM objective and the objective of “satisfy human preferences”!

92

Can we explicitly attempt to **satisfy human preferences?**

93

Reinforcement learning to the rescue

For each LM sample s , imagine we had a way to obtain a human reward of that summary: $R(s) \in \mathbb{R}$, higher is better.

SAN FRANCISCO,
California (CNN) --
A magnitude 4.2
earthquake shook the
San Francisco
...
overturn unstable
objects.

An earthquake hit
San Francisco.
There was minor
property damage,
but no injuries.

$$s_1 \\ R(s_1) = 8.0$$

The Bay Area has
good weather but is
prone to
earthquakes and
wildfires.

$$s_2 \\ R(s_2) = 1.2$$

We want to maximize the expected reward of samples from our LM

94

Optimizing for human preferences

$$\mathbb{E}_{\hat{s} \sim p_\theta(s)} [R(\hat{s})]$$

95

How do we model human preferences?

OpenAI

Model their preferences as a separate (NLP) problem!

An earthquake hit
San Francisco.
There was minor
property damage,
but no injuries.

s_1
 $R(s_1) = 8.0$ 

The Bay Area has
good weather but is
prone to
earthquakes and
wildfires.

s_2
 $R(s_2) = 1.2$ 

Train an LM $RM_\phi(s)$ to
predict human
preferences from an
annotated dataset, then
optimize for RM_ϕ instead.

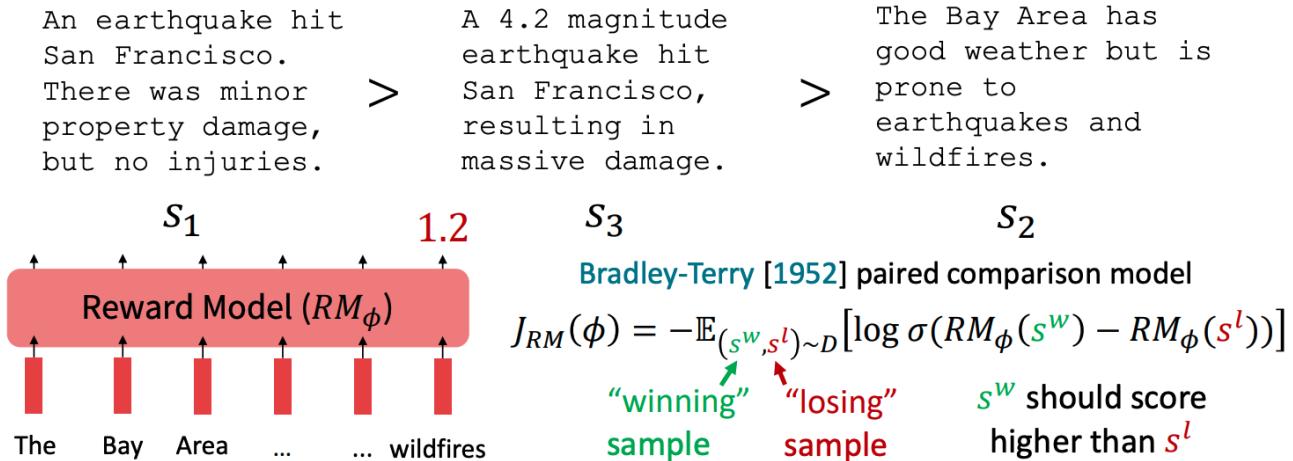
How do we model human preferences?

A 4.2 magnitude
earthquake hit
San Francisco,
resulting in
massive damage.

s_3
 $R(s_3) = 4.1? \quad 6.6? \quad 3.2?$

How do we model human preferences?

Asking pairwise comparisons



RLHF: Putting it all together [Christiano et al., 2017;
Stiennon et al., 2020]

From GPT3 to InstructGPT



PROMPT *Explain the moon landing to a 6 year old in a few sentences.*

COMPLETION GPT-3

Explain the theory of gravity to a 6 year old.

Explain the theory of relativity to a 6 year old in a few sentences.

Explain the big bang theory to a 6 year old.

Explain evolution to a 6 year old.

InstructGPT

People went to the moon, and they took pictures of what they saw, and sent them back to the earth so we could all see them.

100

ChatGPT: Instruction
finetuning + RLHF for
dialog agents



101

From:

Tehran is located in _____

To:

ChatGPT		
 Examples	 Capabilities	 Limitations
"Explain quantum computing in simple terms"	Remembers what user said earlier in the conversation	May occasionally generate incorrect information
"Got any creative ideas for a 10 year old's birthday?"	Allows user to provide follow-up corrections	May occasionally produce harmful instructions or biased content
"How do I make an HTTP request in Javascript?"	Trained to decline inappropriate requests	Limited knowledge of world and events after 2021

102

Practical Applications

Use in Education

Industry Applications

Accessibility and Language Support

103

Limitations of RL + Reward Modeling

Human preferences are unreliable!

"Reward hacking" is a common problem in RL

Google shares drop \$100 billion after its new AI chatbot makes a mistake

February 9, 2023 · 10:15 AM ET

<https://www.npr.org/2023/02/09/1155650909/google-chatbot--error-bard-shares>

Bing AI hallucinates the Super Bowl



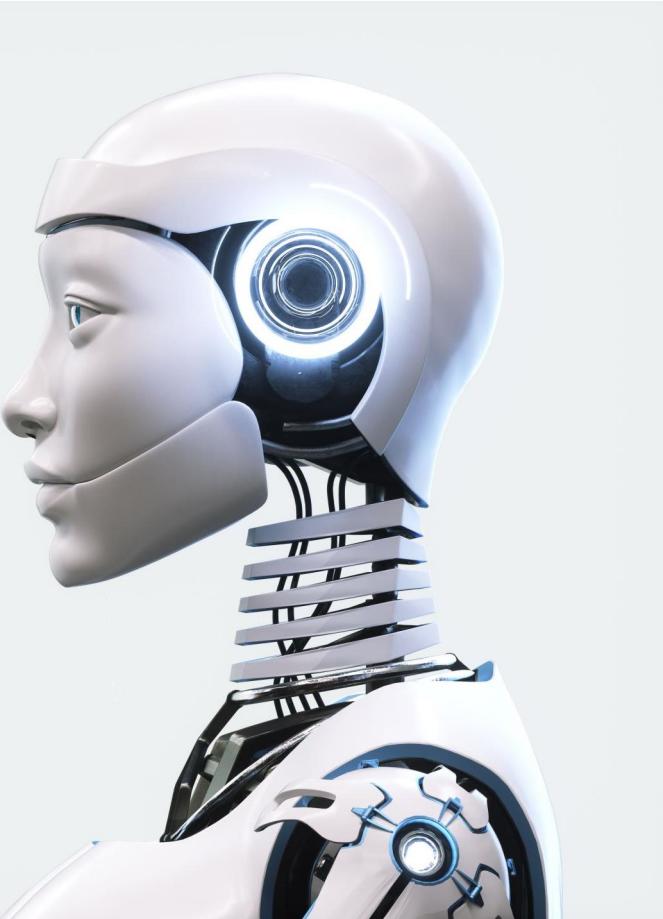
<https://news.ycombinator.com/item?id=34776508>

<https://apnews.com/article/kansas-city-chiefs-philadelphia-eagles-technology-science-82bc20f207e3e4cf81abc6a5d9e6b23a>

104

Chatbots are rewarded to produce responses that seem **authoritative** and **helpful**, regardless of truth

Models of human preferences are even more unreliable!



105



ChatGPT hallucinates

Prompt	Where do people smoke more on average than in Russia?
ChatGPT Generation	<p>There are several countries where people smoke more on average than in Russia. According to data from the World Health Organization, the following countries have higher smoking rates than Russia:</p> <p>Indonesia: 55.2% of adults smoke China: 52.9% of adults smoke Serbia: 47.5% of adults smoke Serbia and Montenegro: 45.6% of adults smoke Montenegro: 45.6% of adults smoke Greece: 44.5% of adults smoke Bosnia and Herzegovina: 41.8% of adults smoke Kazakhstan: 41.3% of adults smoke Kosovo: 41.3% of adults smoke Moldova: 40.9% of adults smoke</p> <p>It's important to note that these figures are estimates and may vary slightly from other sources.</p>
Hallucination	Extrinsic - the generated answer includes wrong statistics. (Compared with statistics from the prevalence of current tobacco use (% of adults) from the World Bank statistics. ¹⁶⁾

Numbers marked in red are wrong, ChatGPT imagines them.

ChatGPT has a timestamp

AS

Who won the most recent FIFA World Cup?

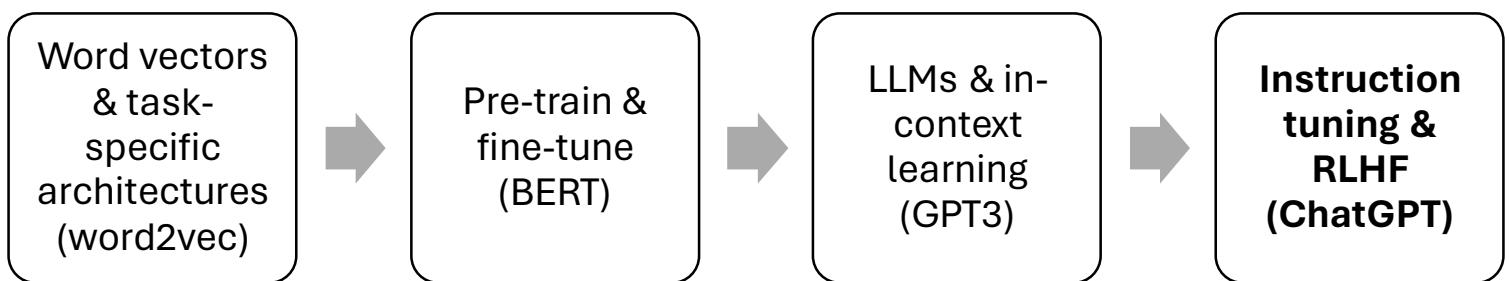


The most recent FIFA World Cup was held in 2018, and the winner was the French national team. They defeated Croatia 4-2 in the final to win their second World Cup title. The tournament was held in Russia from June 14 to July 15, 2018, and featured 32 teams from around the world.

The answer should have been [Argentina](#), but it is not always trained on the most recent data.

108

Recent history



109