

Generative AI: Foundations

CPEN 322

Agenda

00

FOUNDATIONS
OF ML

Speedrun
key ideas in ML

01

TRANSFORMER
ARCHITECTURE

Core ideas and
notable examples

02

NOTABLE LLMs

T5, GPT, Chinchilla, et al

03

TRAINING
& INFERENCE

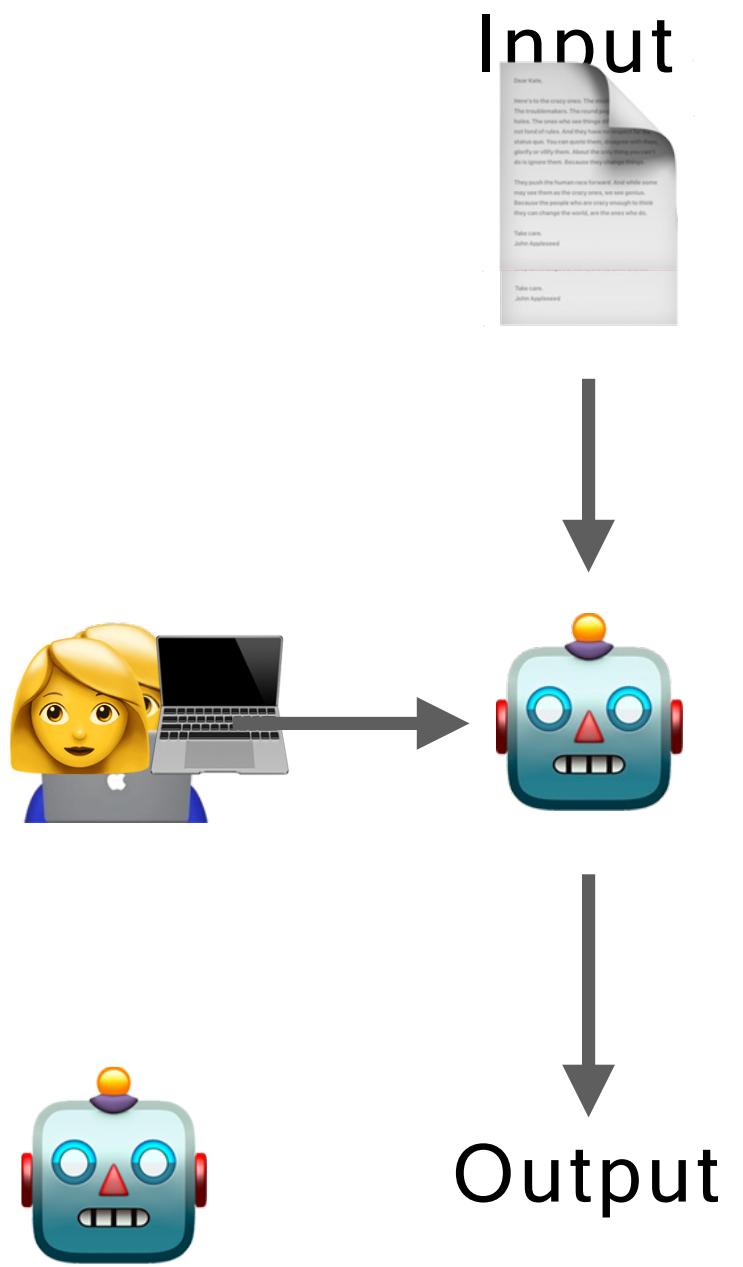
Running a Transformer

00

Foundations of Machine Learning

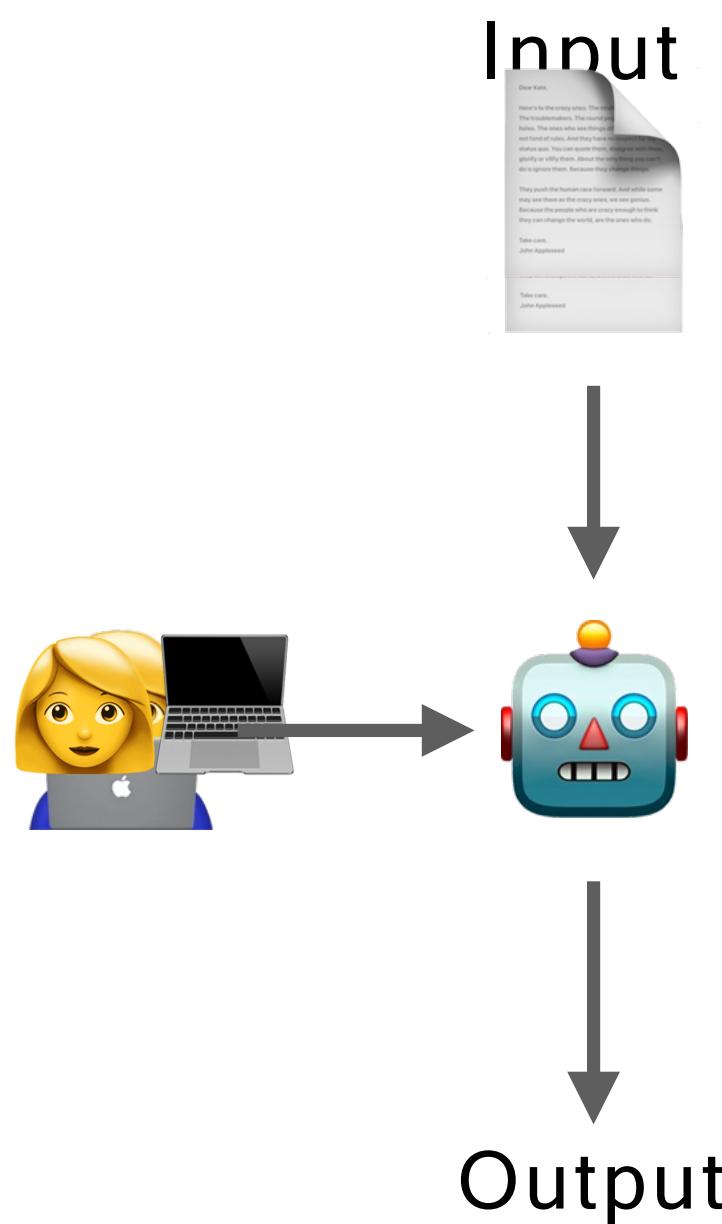


Traditional Programming

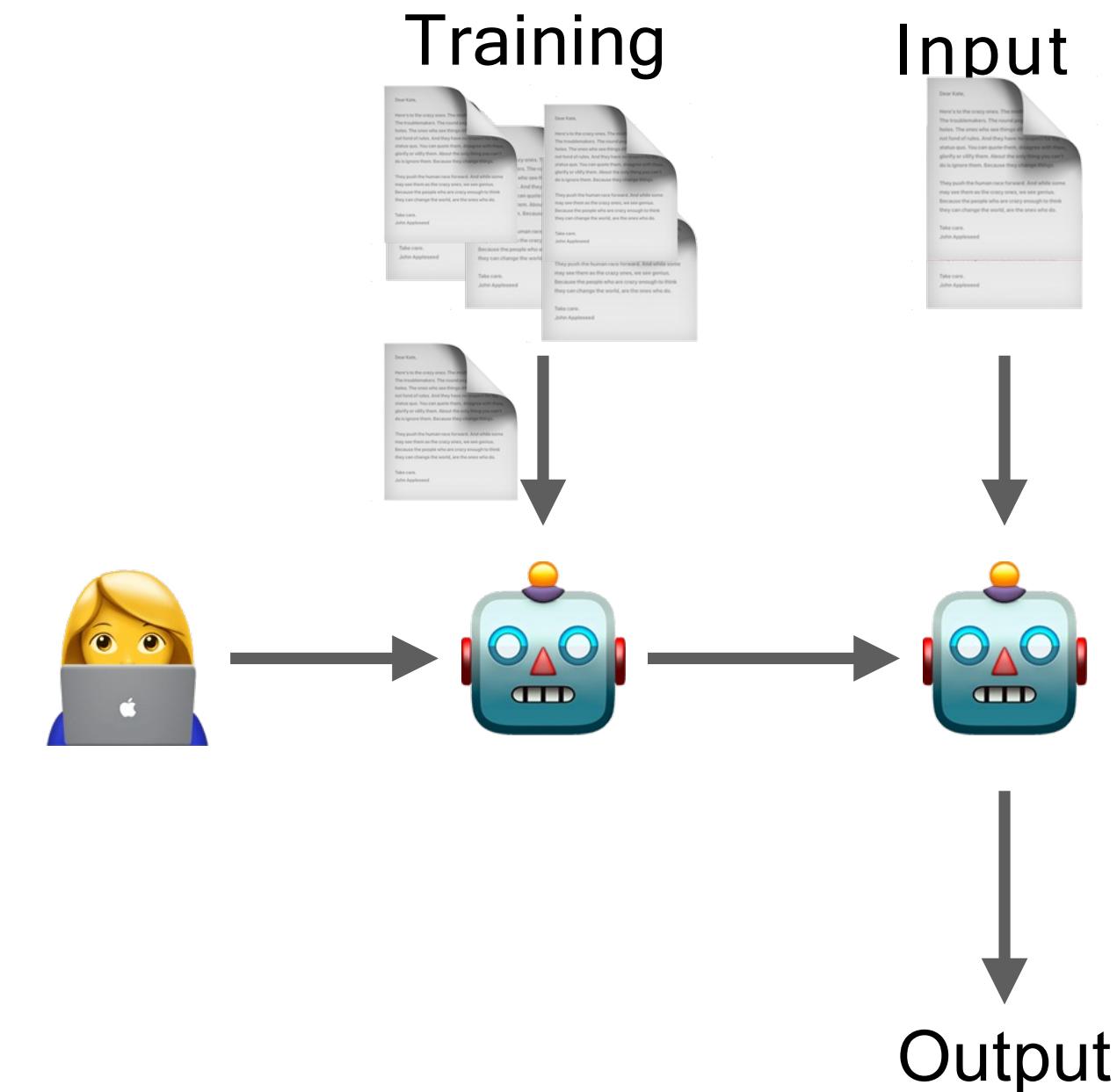


Software 1.0

Traditional Programming vs Machine Learning



Software 1.0



Software 2.0

Types of Machine Learning

Unsupervised Learning

Learn structure of data to predict (e.g., clustering)

"This product does what it is supposed "

Supervised Learning

Learn how data maps to labels to recognize or predict



→ *cat*



→ "Hey
Siri"

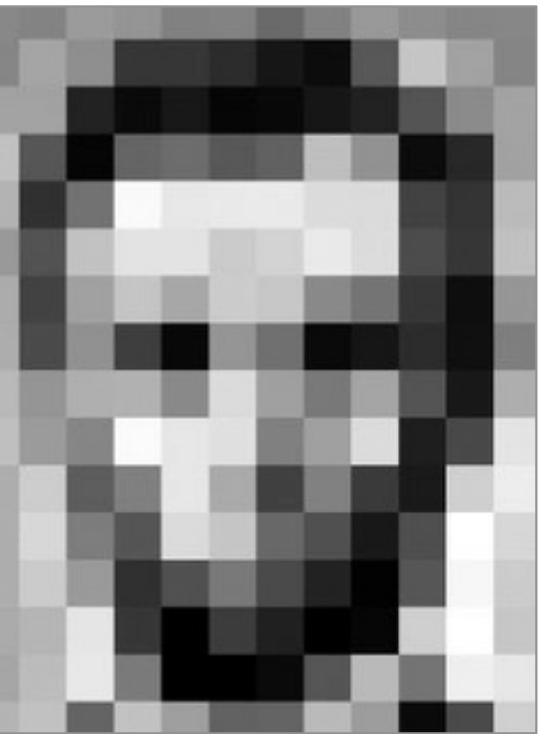
Reinforcement Learning

Learn how to act in an environment to obtain reward



Embedding: Inputs and outputs are always numbers

What we see



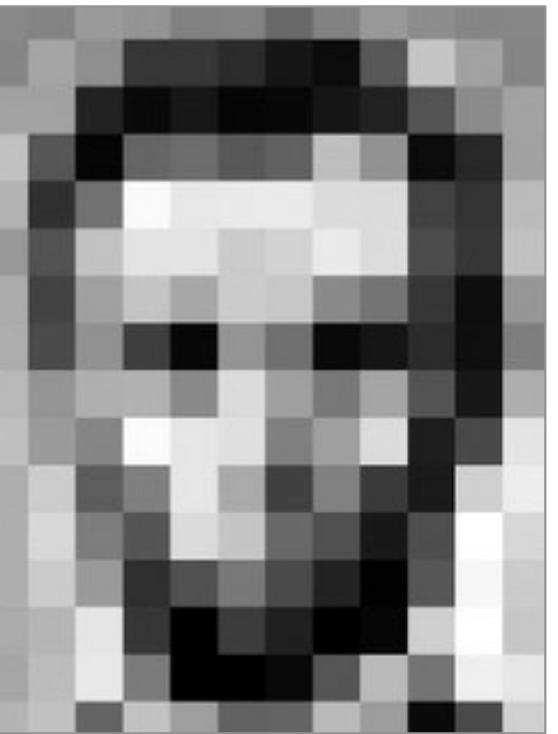
Input

Output

"Lincoln"

Embedding: Inputs and outputs are always numbers

Input



What we see

Output

"Lincoln"

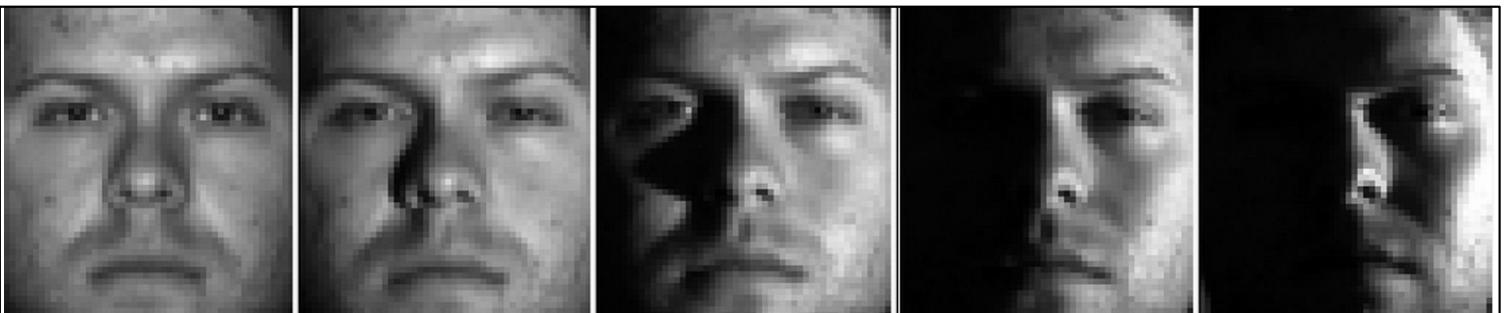
157	153	174	168	150	152	129	161	172	161	156
156	182	163	74	75	62	33	17	110	210	180
180	180	50	14	34	6	10	33	48	106	159
206	109	5	124	131	111	120	204	166	15	56
194	68	137	251	237	239	239	228	227	87	71
172	105	207	233	233	214	220	239	228	98	74
188	88	179	209	185	215	211	198	139	75	20
189	97	165	84	10	168	134	11	31	62	22
199	168	191	193	158	227	178	143	182	106	36
205	174	195	252	236	231	149	178	228	43	95
190	216	116	149	236	187	86	150	79	38	218
190	224	147	108	227	210	127	102	36	101	255

What the machine "sees"

[76, 105, 110, 99, 111, 108, 110]

Why is this hard?

- Infinite variety of inputs can all mean the same thing
- Meaningful differences can be tiny
- Structure of the world is complex



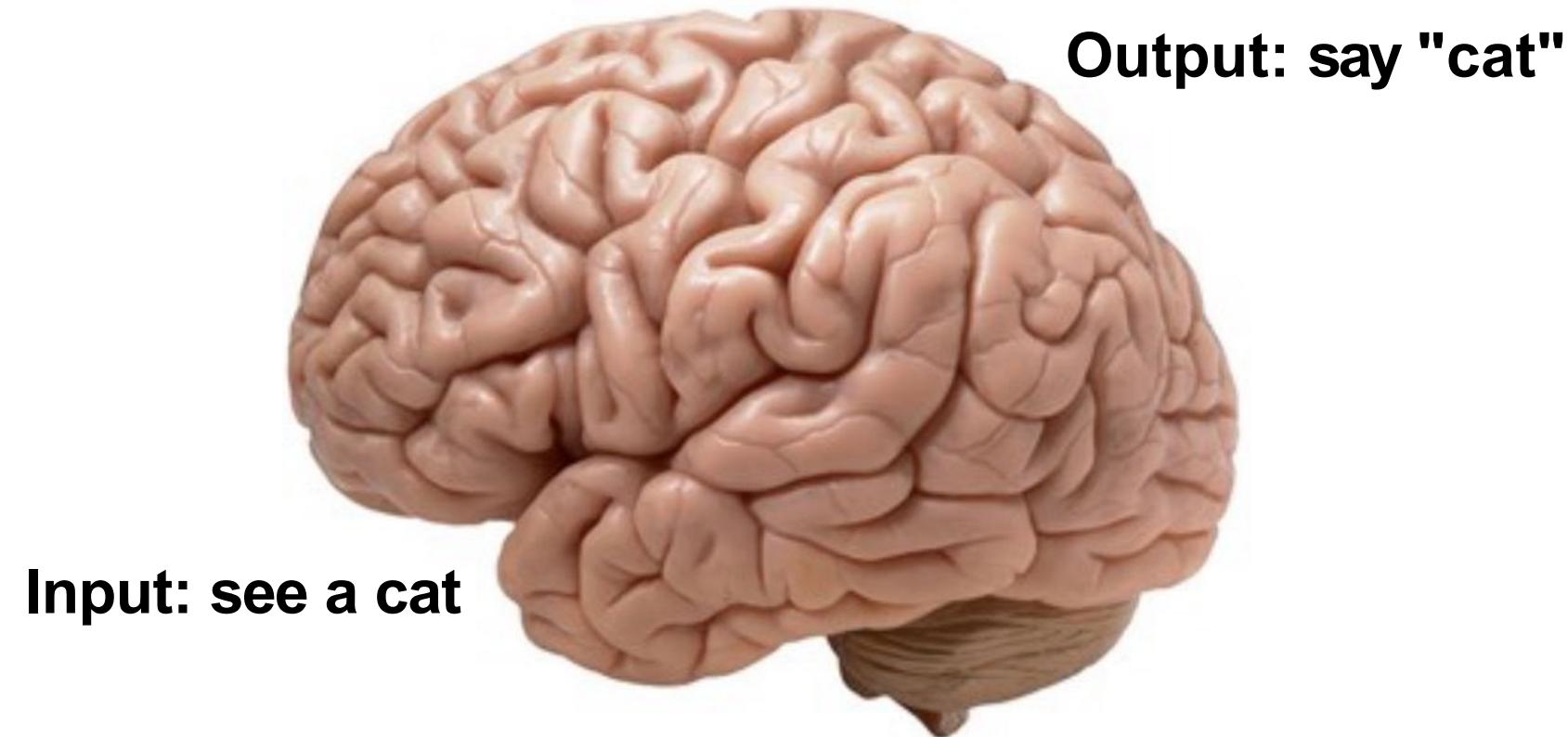
"I loved this movie"

"As good as The Godfather"

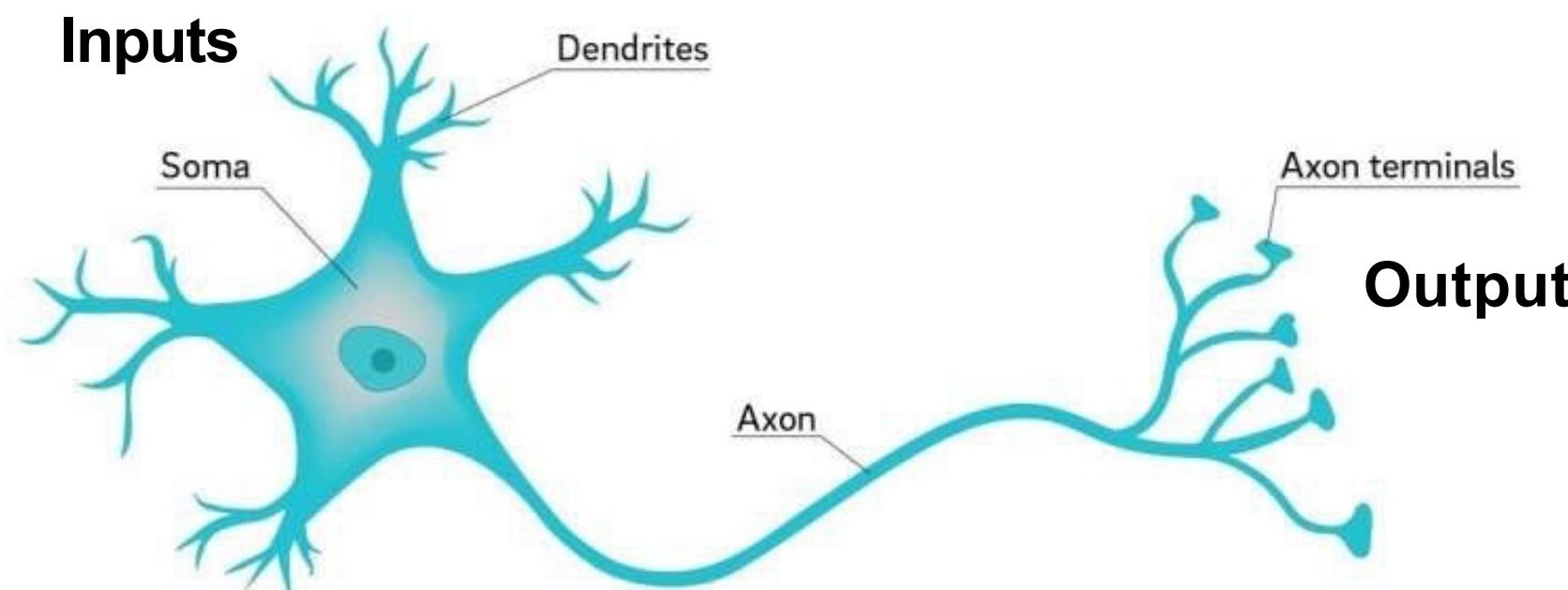
How is it done?

- Many methods for Machine Learning
 - Logistic Regression
 - Support Vector Machines
 - Decision Trees
- But one is dominant
 - Neural Networks (also called Deep Learning)

Inspiration



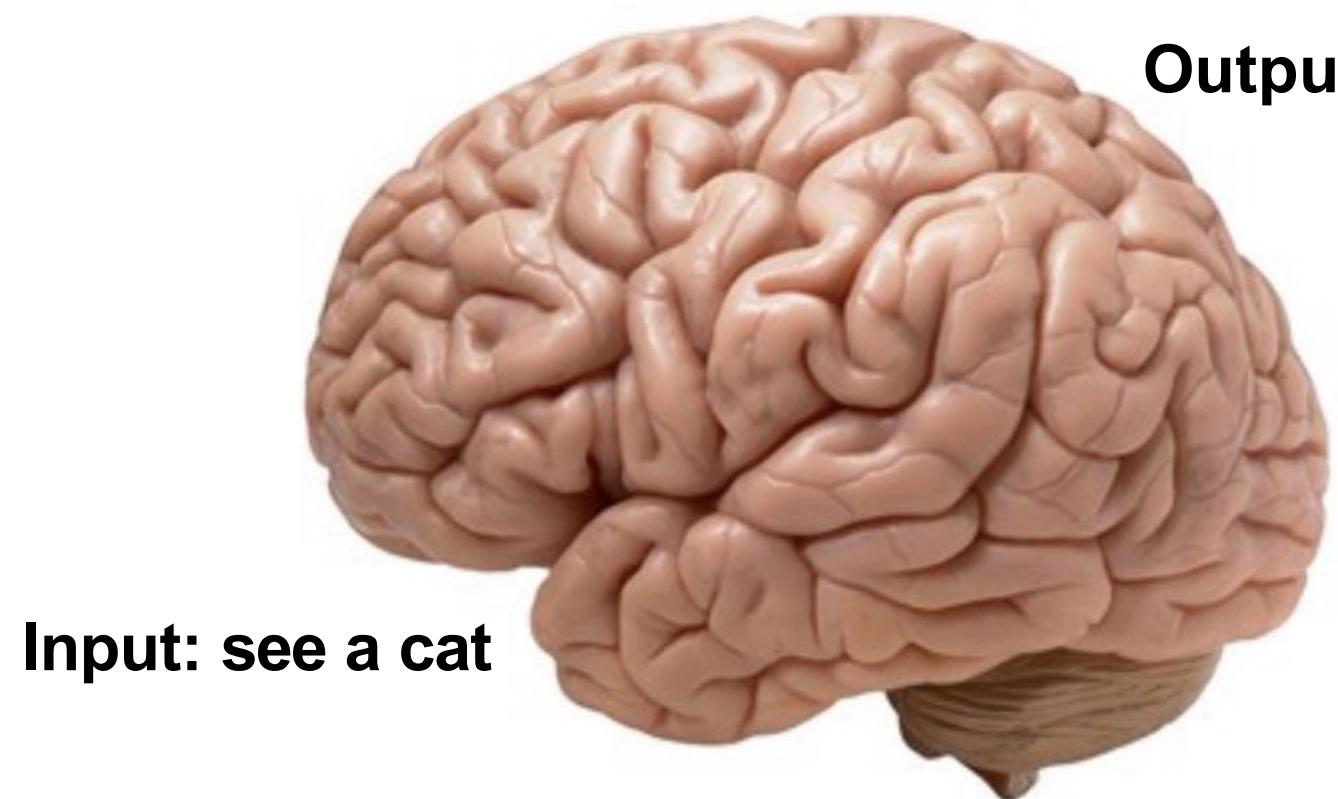
<https://www.the-scientist.com/the-nutshell/what-made-human-brains-so-big-36663>



<https://medicalxpress.com/news/2018-07-neuron-axons-spindly-theyre-optimizing.html>

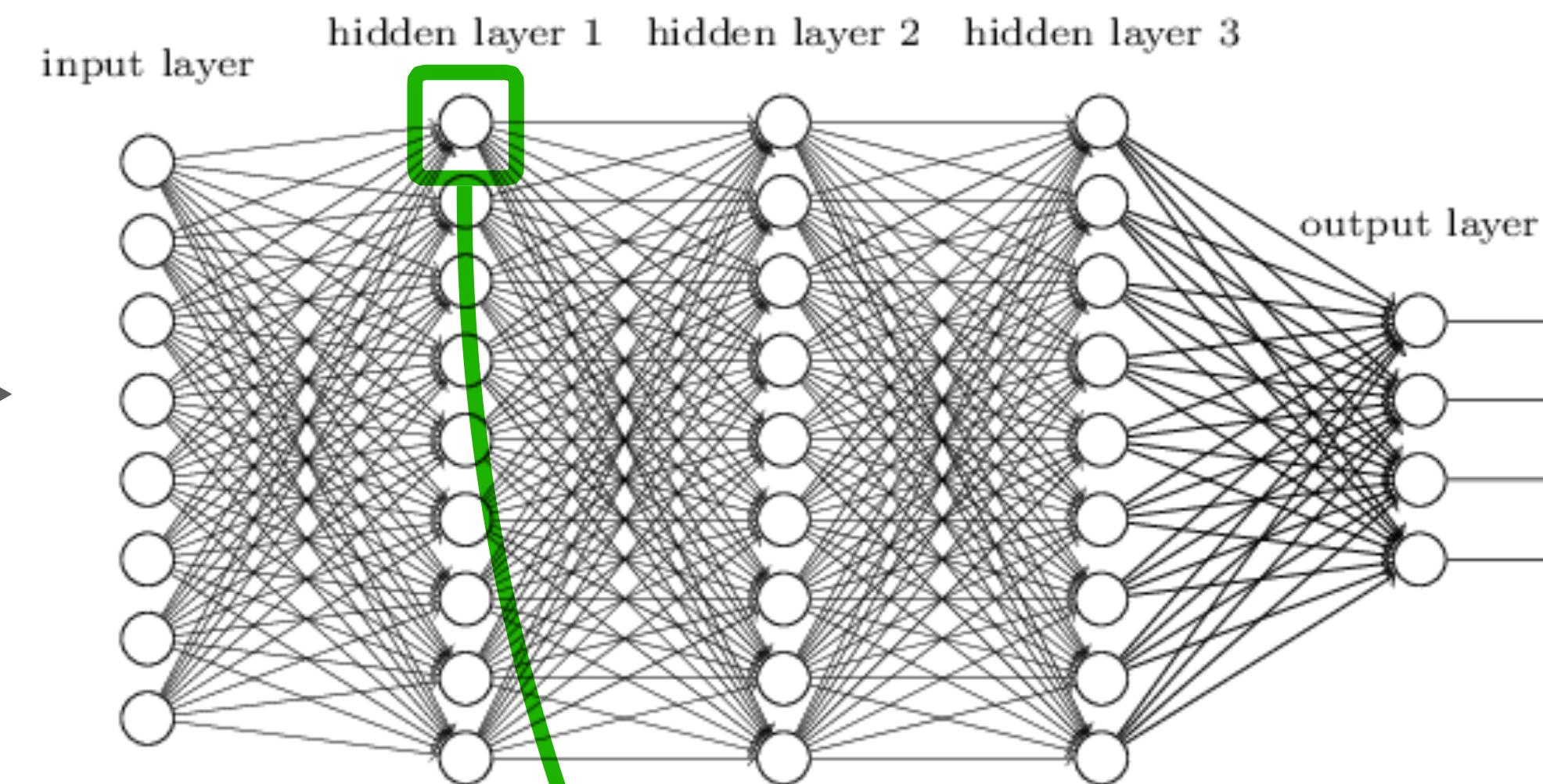
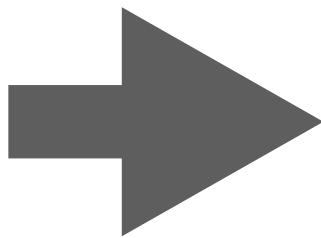
- Inspired by what we know to be intelligent: the **brain**
- The brain is composed of billions of **neurons**
- Each neuron receives electrical **inputs** and sends an electrical **output**
- The brain itself has high-level inputs and outputs

Formalization

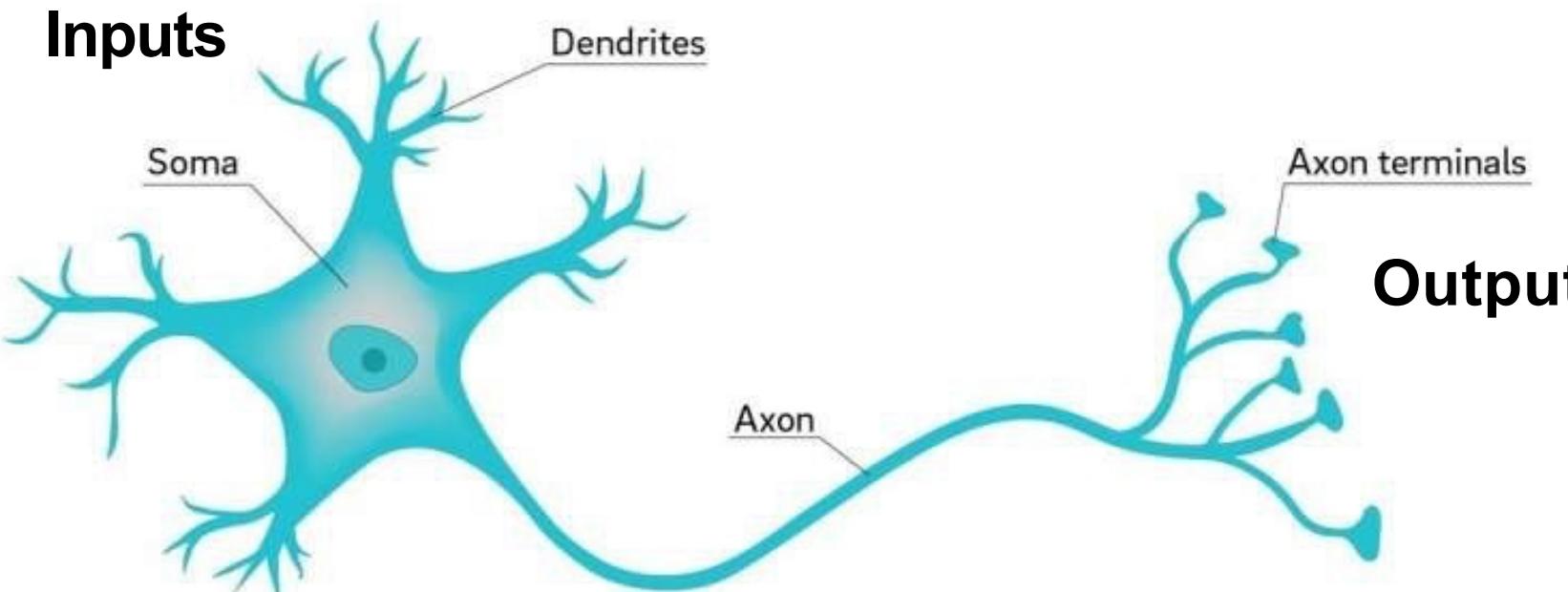


Input: see a cat

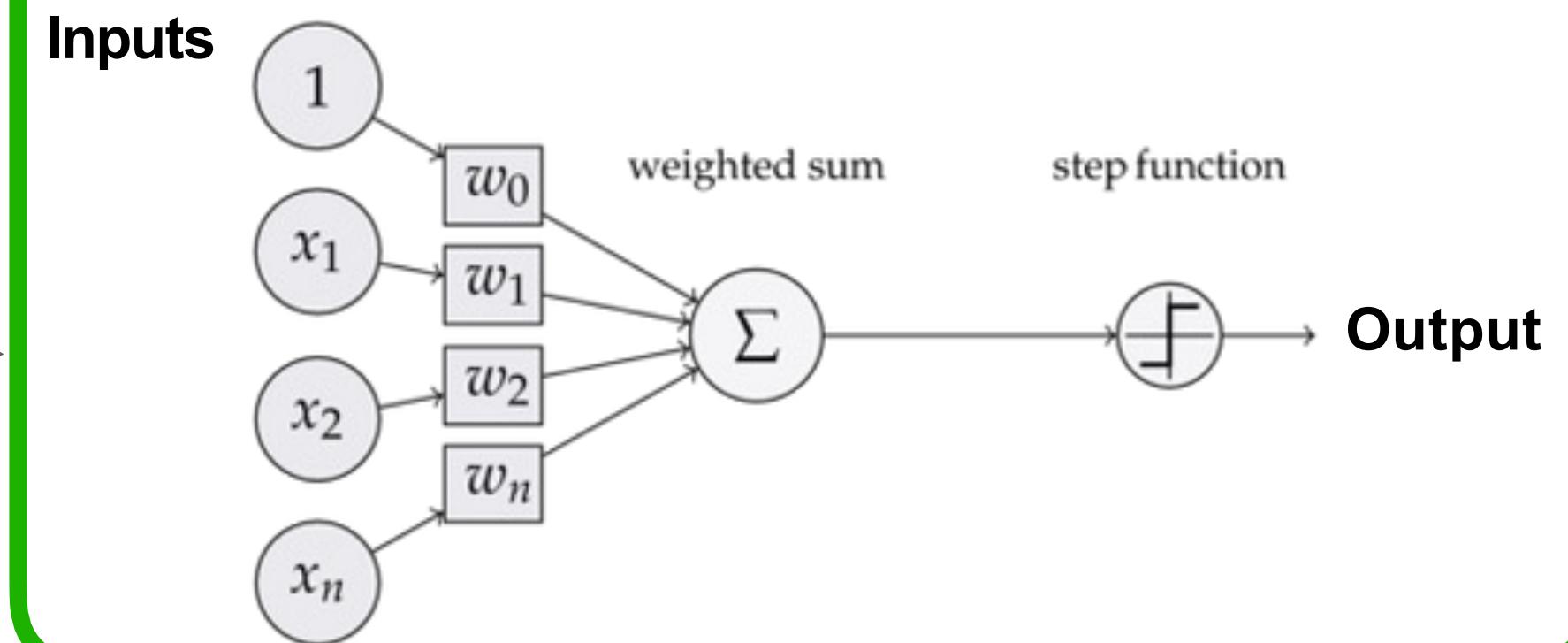
Output: say "cat"



<https://www.the-scientist.com/the-nutshell/what-made-human-brains-so-big-36663>

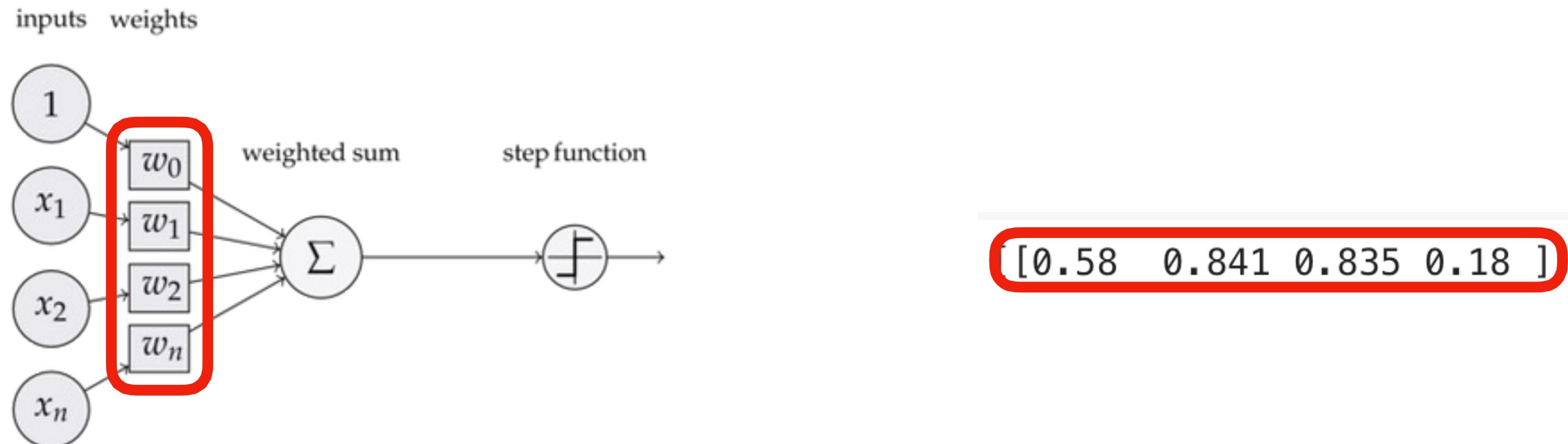


<https://medicalxpress.com/news/2018-07-neuron-axons-spindly-theyre-optimizing.html>



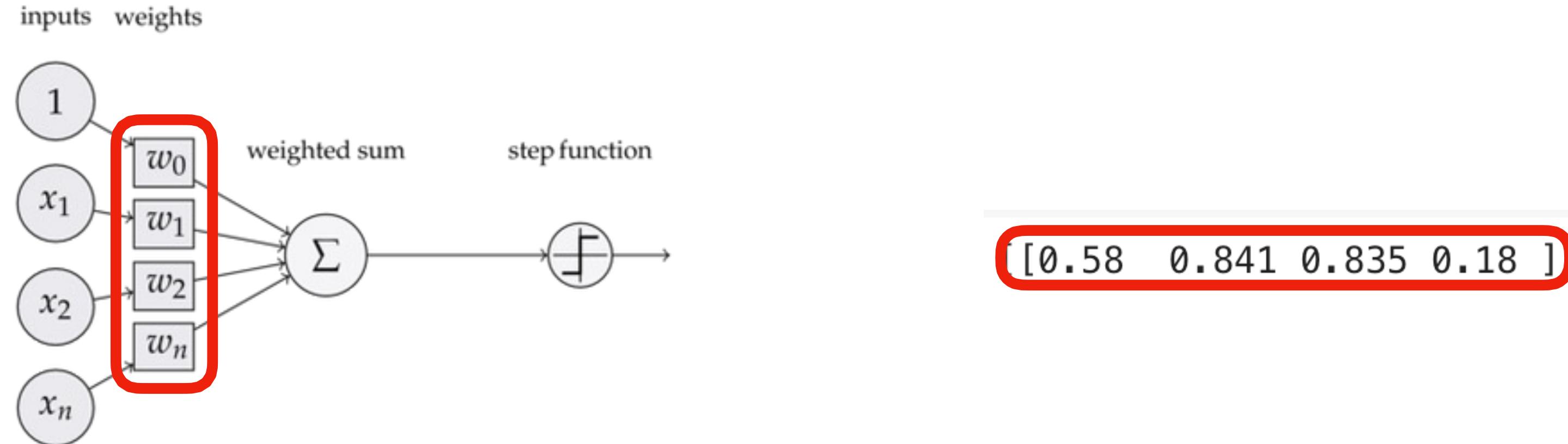
<https://www.jessicayung.com/explaining-tensorflow-code-for-a-multilayer-perceptron/>

A "perceptron" is a vector of numbers



<https://www.jessicayung.com/explaining-tensorflow-code-for-a-multilayer-perceptron/>

Step Function

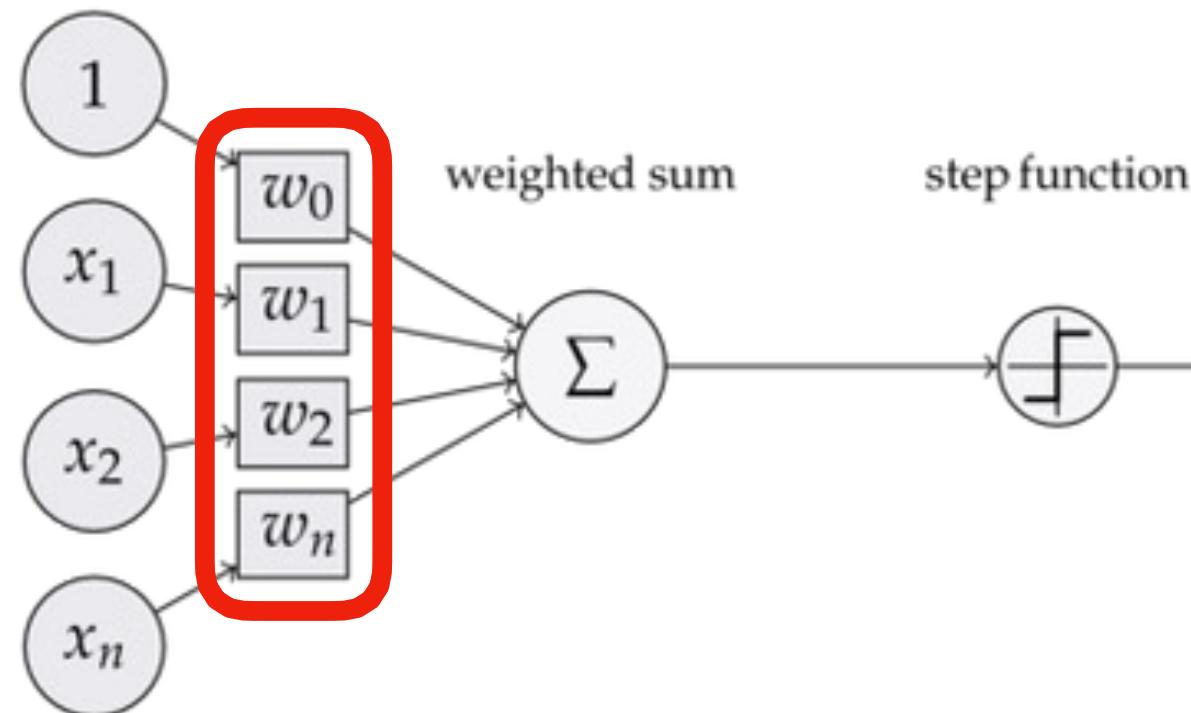


<https://www.jessicayung.com/explaining-tensorflow-code-for-a-multilayer-perceptron/>

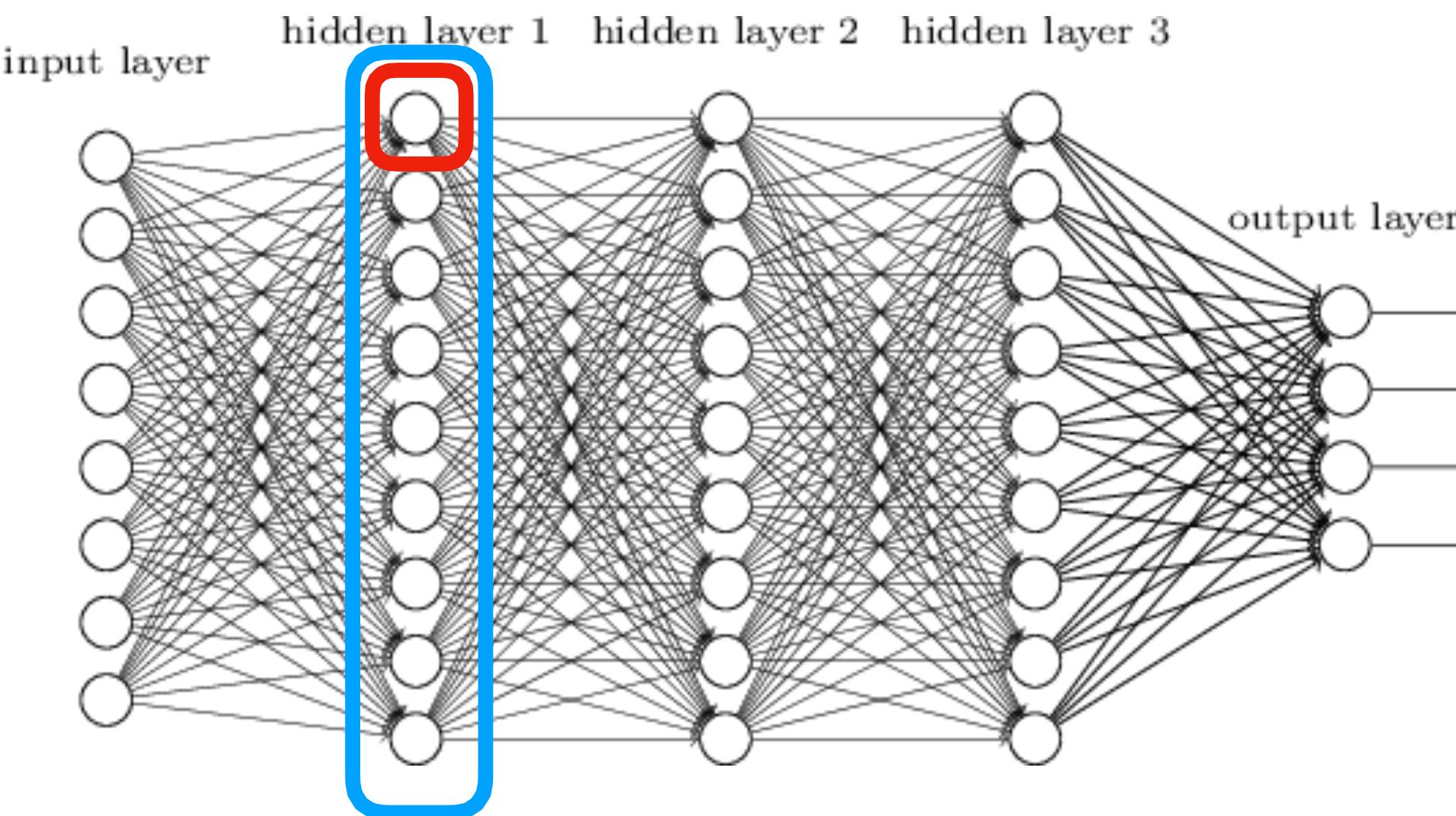
- Step function: also called threshold function
 - Decides based on the weighted sum if the output is 0 or 1
- In modern deep learning, the step function is less often used as it is not differentiable at the threshold point
- Instead, smoother functions that approximate the step function but retain differentiability are used
- For example: **sigmoid** outputs a continuous value in the **range (0, 1)**

A "layer" is a matrix of numbers

inputs weights



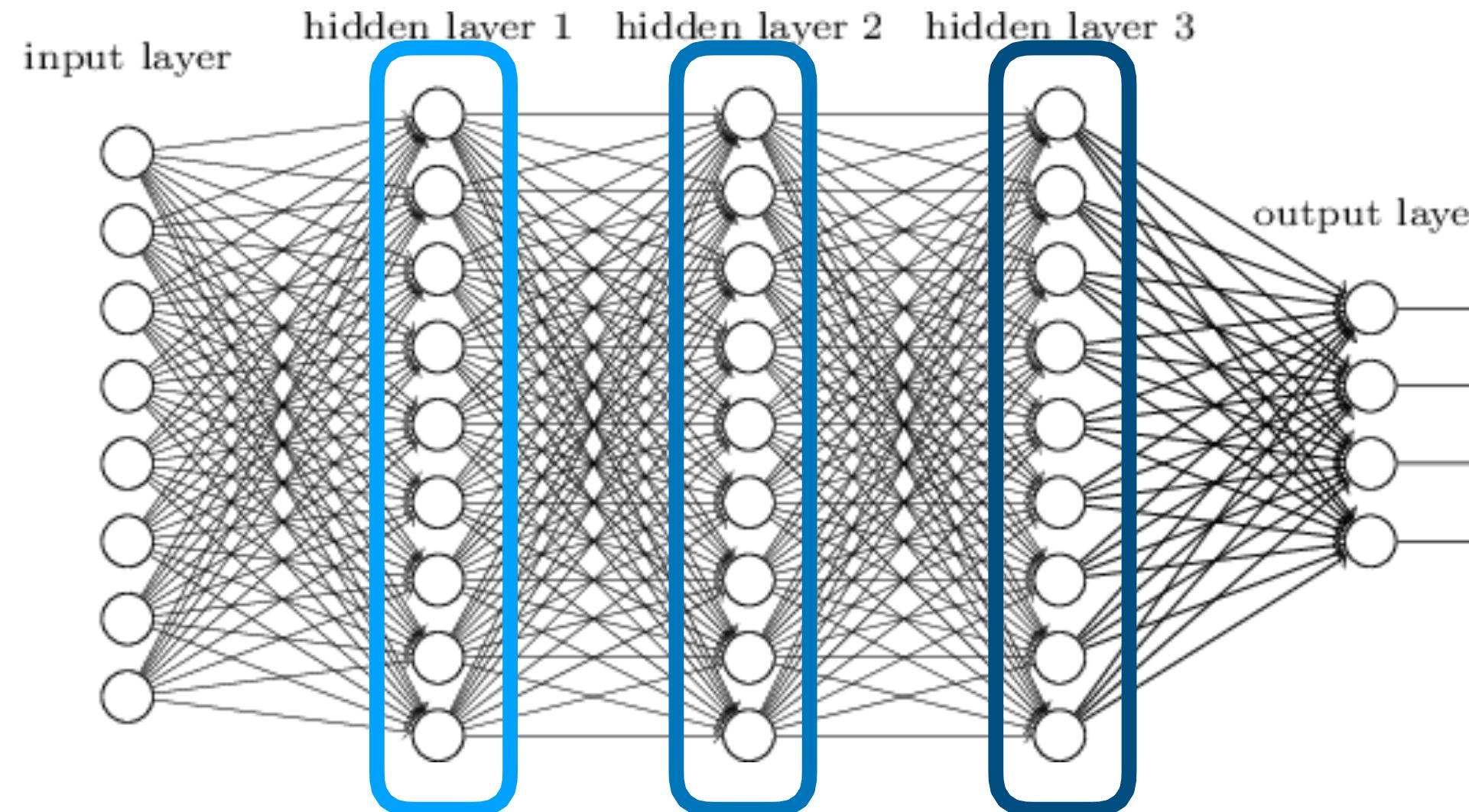
<https://www.jessicayung.com/explaining-tensorflow-code-for-a-multilayer-perceptron/>



[0.58 0.841 0.835 0.18]
[0.405 0.813 0.309 0.562]
[0.422 0.229 0.46 0.152]
[0.673 0.429 0.441 0.243]
[0.9 0.744 0.234 0.856]
[0.971 0.486 0.175 0.248]
[0.258 0.588 0.478 0.266]
[0.236 0.496 0.077 0.557]
[0.413 0.322 0.372 0.741]]

The neural network is a set of matrices

Called "parameters" or "weights"

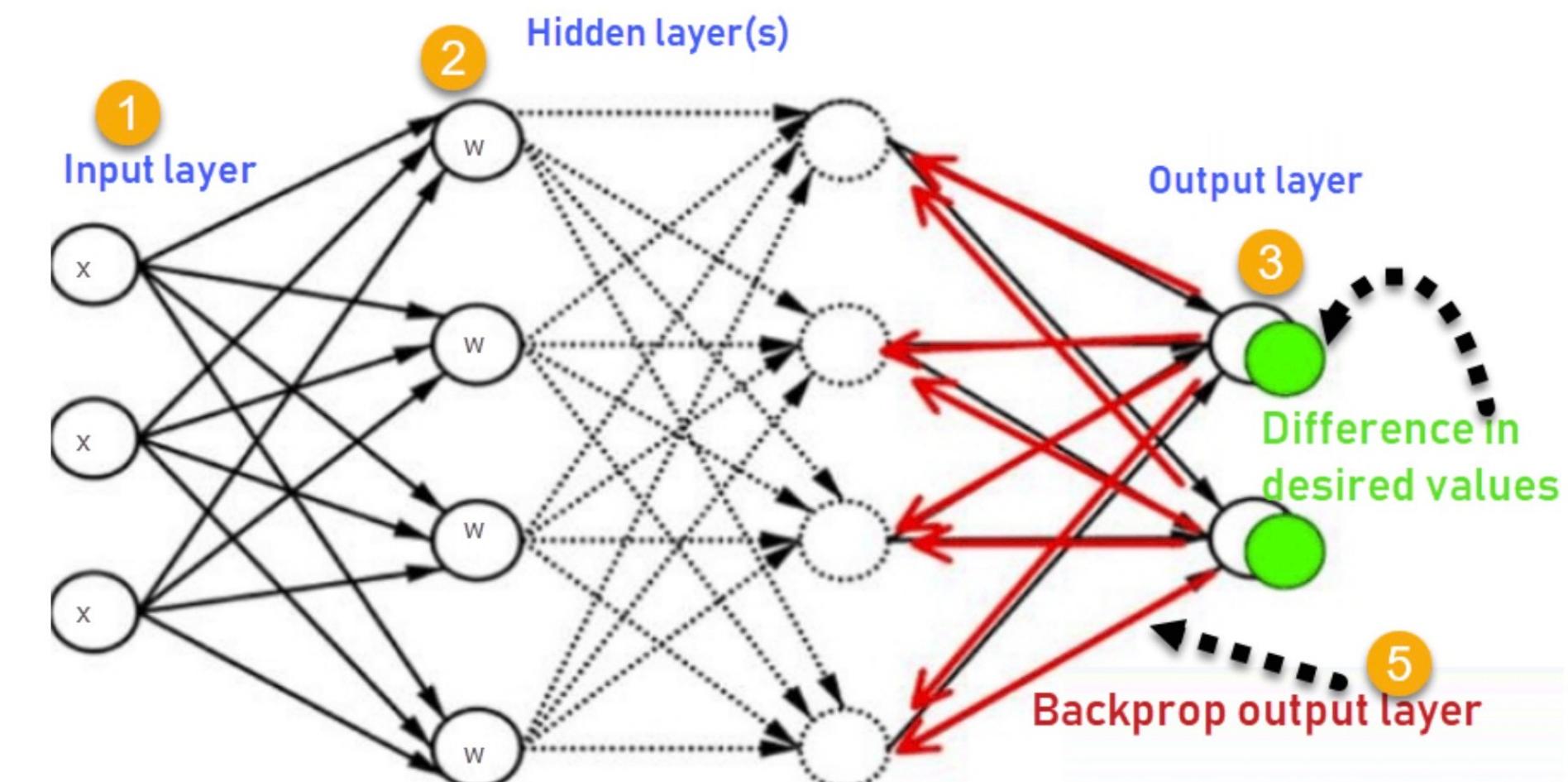


[[0.32 0.866 0.151 0.555]
[0 0.58 0.841 0.835 0 18 1
[0 [[0.32 0.866 0.151 0.555]
[0 [0.467 0.414 0.562 0.555]
[0 [0.186 0.937 0.131 0.448]
[0 [0.892 0.567 0.111 0.678]
[0 [0.483 0.478 0.424 0.844]
[0 [0.774 0.963 0.213 0.569]
[0 [0.014 0.509 0.21 0.26]
[0 [0.142 0.991 0.105 0.211]
[0 0.789 0.242 0.774 0.019]]

NN operations are just matrix multiplications.
GPUs are really fast at matrix multiplications.

Training

- Dataset contains n items
- Batch: a subset of n to use for training
 - If n is large, it won't fit in memory
- Epoch: going over the entire n items in the dataset

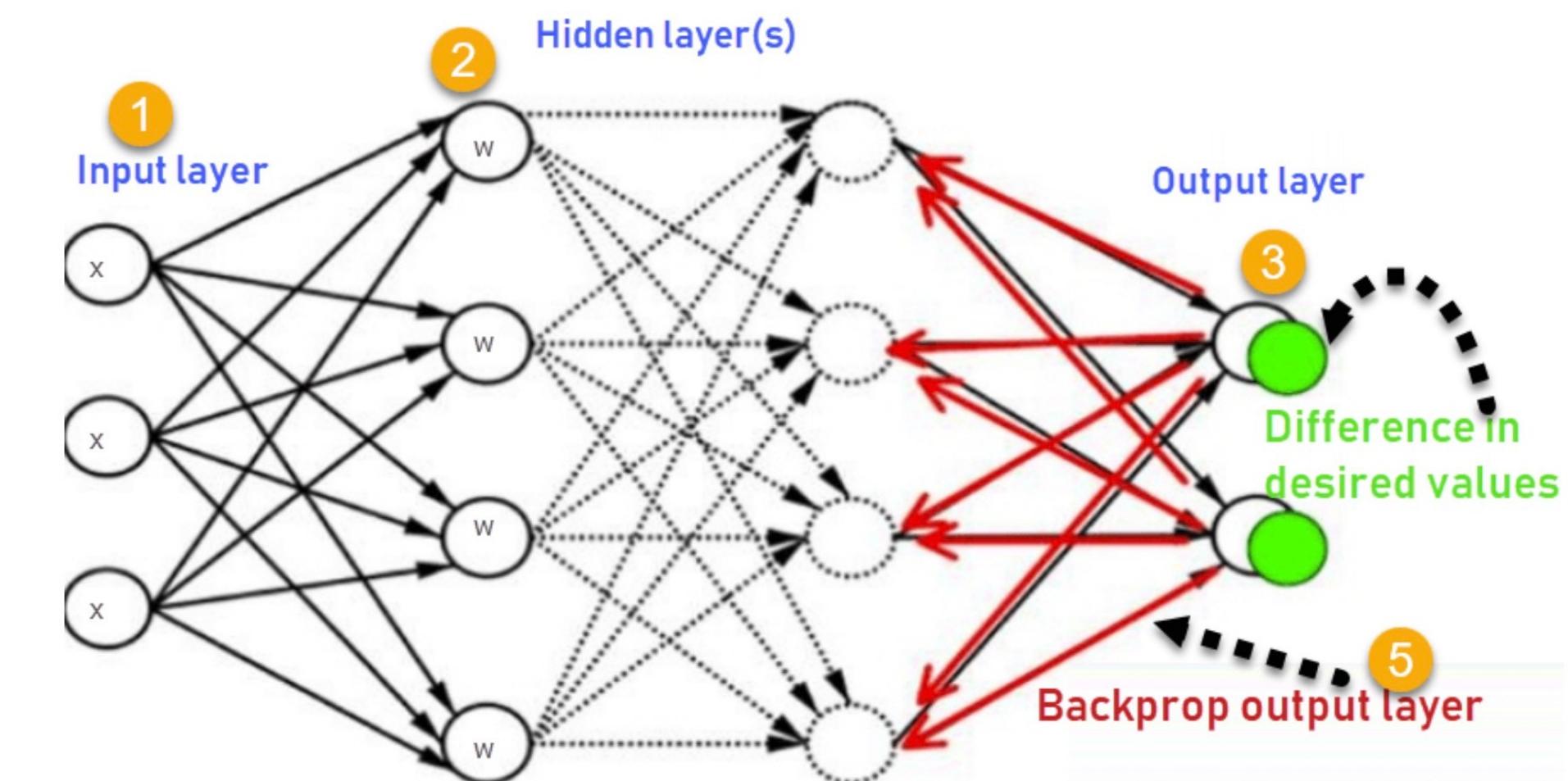
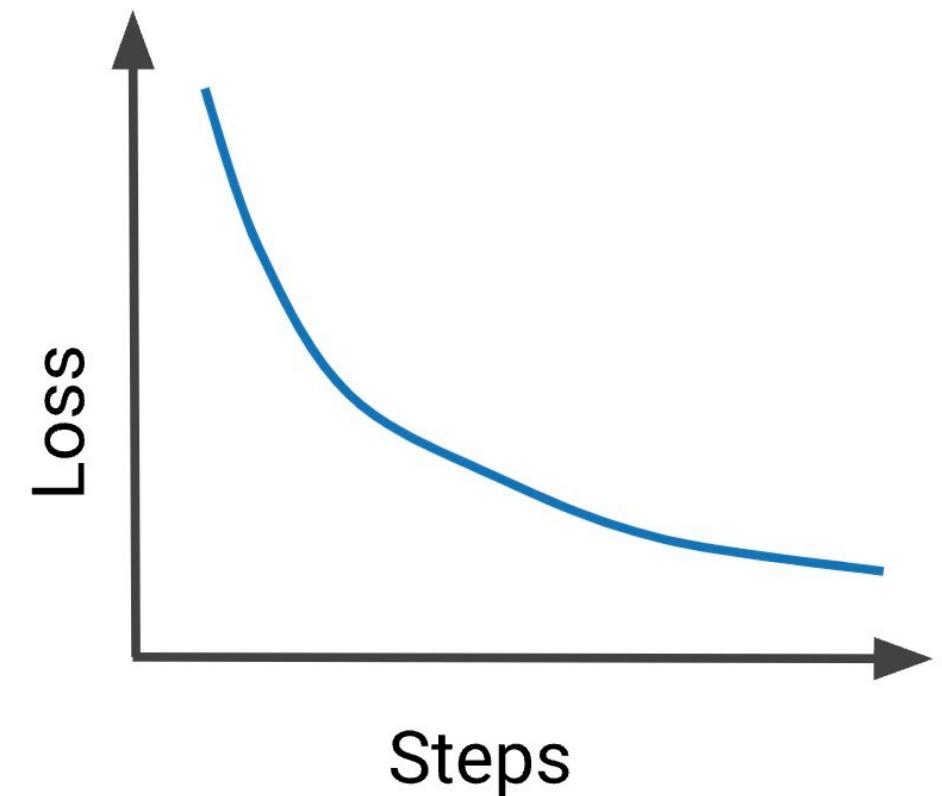


<https://www.guru99.com/backpropogation-neural-network.html>

Training

- Data X (e.g. images), labels y (e.g. labels)
- Take a little batch of data x:
 - Use the current model to make a prediction $x \rightarrow y'$
 - Compute loss(y, y')
 - Back-propagate the loss through all the layers of the model
- Repeat until loss stops decreasing

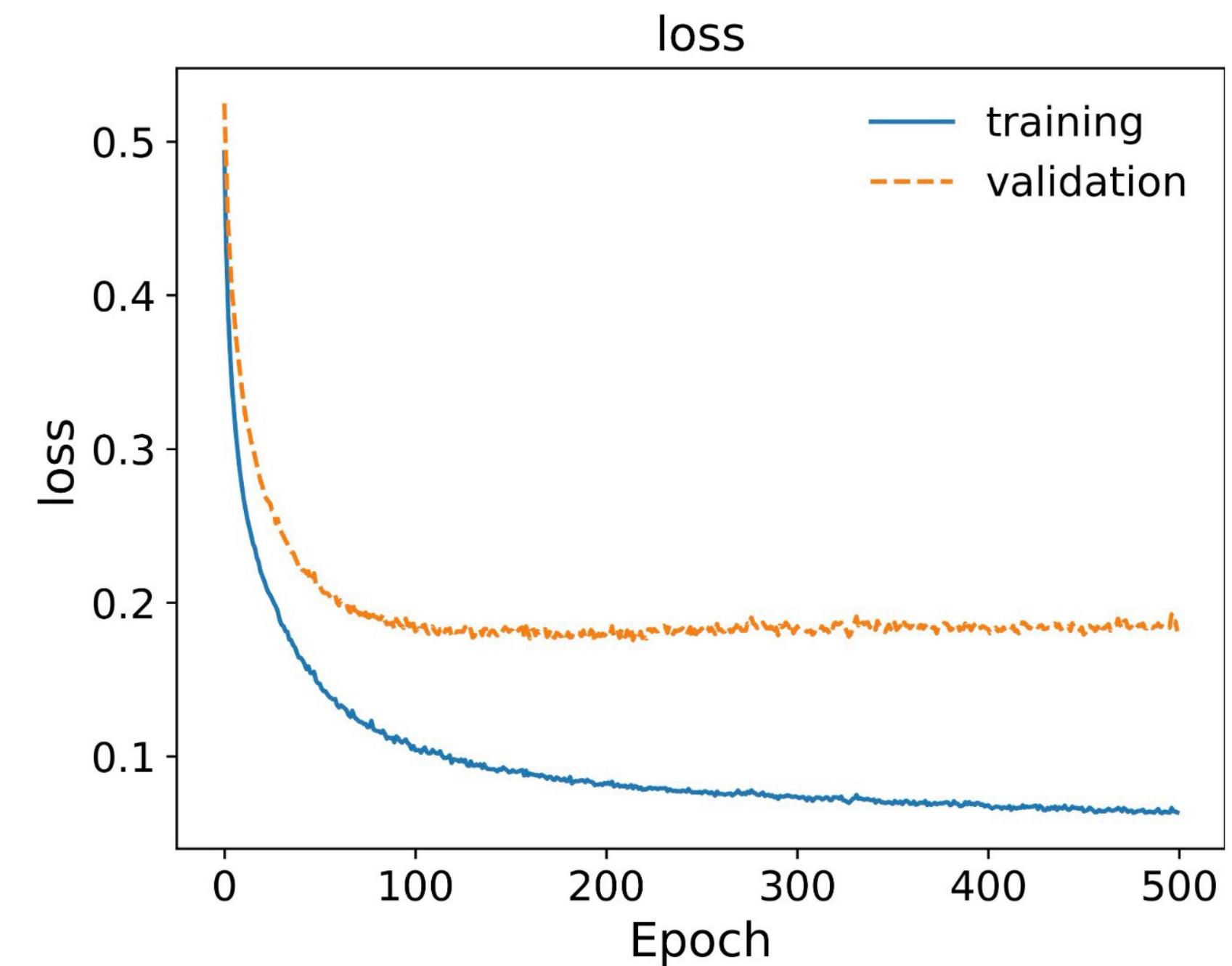
$$loss_{CE} = - \sum (y_i * \log(y'_i))$$



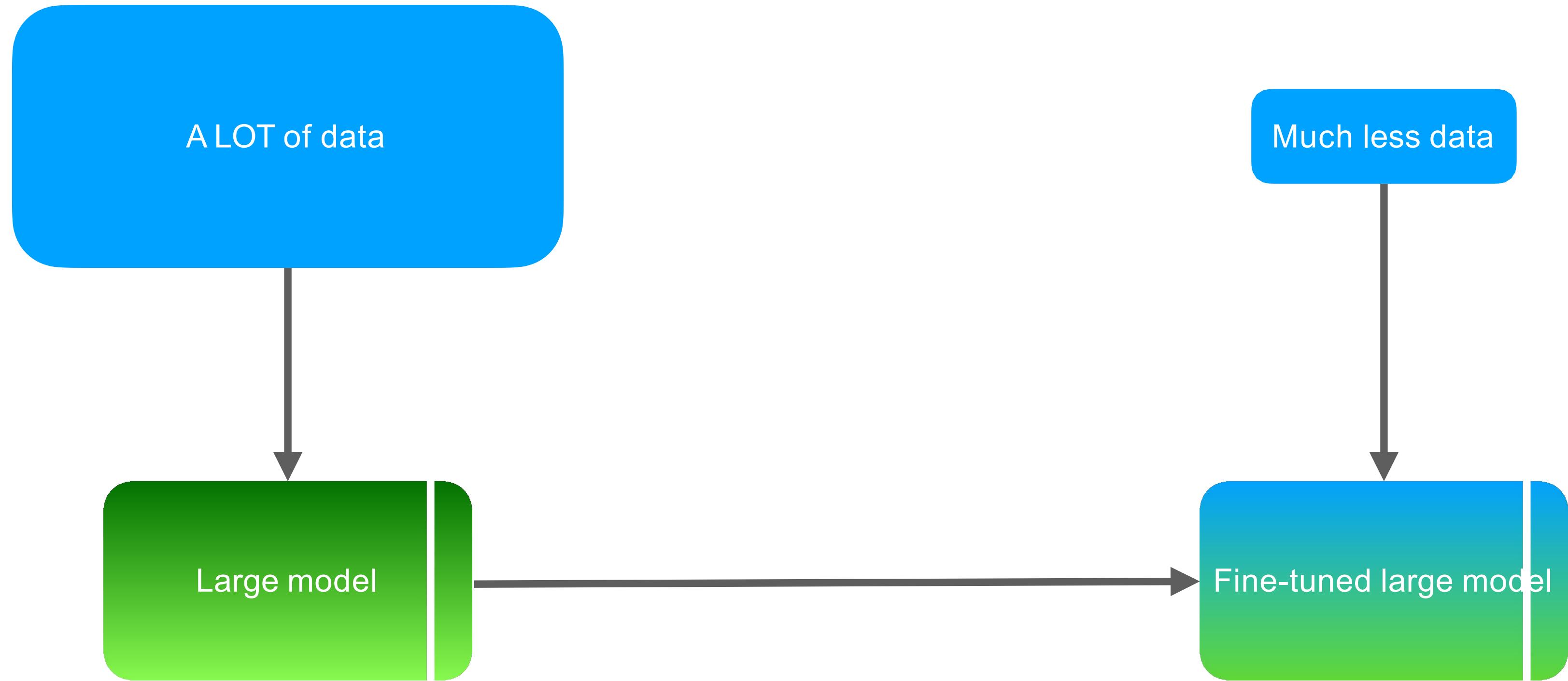
<https://www.guru99.com/backpropogation-neural-network.html>

Dataset Splitting

- Split (X, y) into training ($\sim 80\%$), validation ($\sim 10\%$), and test ($\sim 10\%$) sets
- Validation set is for
 - ensuring that training is not "overfitting"
 - setting hyper-parameters of the model (e.g. number of parameters)
- Test set is for measuring validity of predictions on new data



THIS APPLIES TO YOUR
EXPERIMENTATION
WITH PROMPTS!



Pre-training:
slow training on a lot of data

Fine-tuning:
fast training on a little data

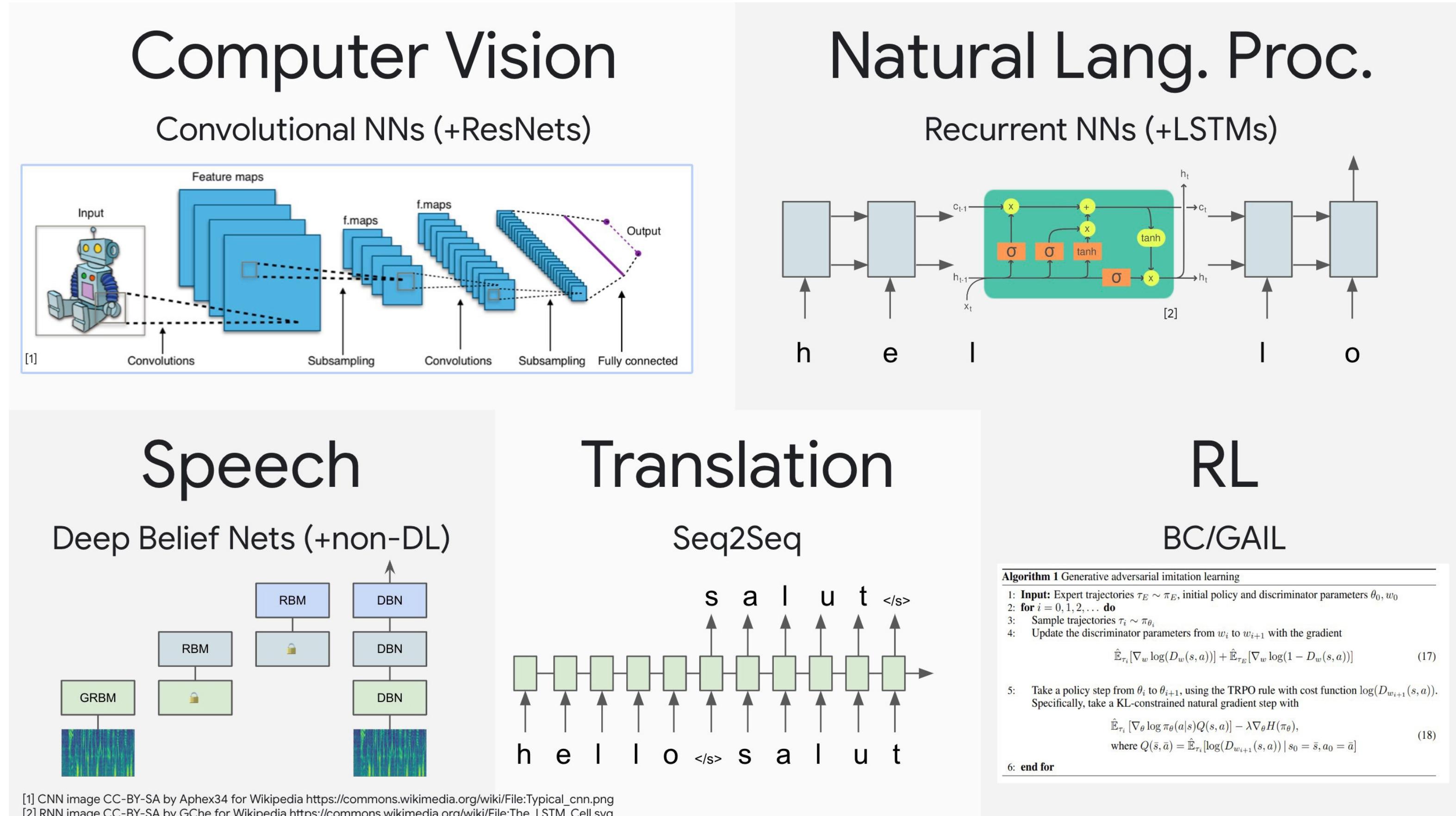
Model Hubs

- People share pre-trained models!
 - 180K models
 - 30K datasets

The screenshot shows the Hugging Face Model Hub interface. At the top, there is a logo, a search bar, and navigation links for Models, Datasets, Spaces, Docs, Solutions, and Pricing. The main area is divided into sections: Tasks (1), Libraries, Datasets, Languages, Licenses, Other, Multimodal, Computer Vision, Natural Language Processing, and Audio. Each section contains sub-links for various model types. On the right, a large list of pre-trained models is displayed, each with a profile picture, name, last updated date, file size, and download count.

Model Name	Last Updated	File Size	Downloads
cardiffnlp/twitter-roberta-base-sentiment	Jan 20	2.39M	155
Seethal/sentiment_analysis_generic_dataset	Apr 18, 2022	1.33M	10
cardiffnlp/twitter-roberta-base-sentiment-latest	Jan 13	1.24M	117
yiyanghkust/finbert-tone	Oct 16, 2022	1.03M	70
ProsusAI/finbert	Oct 2, 2022	699k	199
zhayunduo/roberta-base-stocktwits-finetuned	Apr 18, 2022	347k	10
nlptown/bert-base-multilingual-uncased-sentiment	Apr 18, 2022	260k	108
prithivida/parrot_fluency_model	Jun 24, 2022	210k	
oliverguhr/german-sentiment-bert	Mar 16	170k	19
distilbert-base-uncased-fine-tuned	29 days ago	2.16M	193
cardiffnlp/twitter-xlm-roberta-base	Nov 28, 2022	1.32M	97
papluca/xlm-roberta-base-10langs	Nov 5, 2022	1.08M	77
finiteautomata/bertweet-base	Feb 16	958k	58
j-hartmann/emotion-english	Jan 2	597k	139
prithivida/parrot_adequacy	May 26, 2022	299k	3
roberta-base-openai-detector	13 days ago	237k	68
cross-encoder/ms-marco-MiniLM-L12	Aug 5, 2021	172k	16
yiyanghkust/finbert-esg-9	Oct 16, 2022	155k	11

Before ~2020: each task had its own NN architecture



Now: all is Transformers



Transformer cartoon (DALL-E)

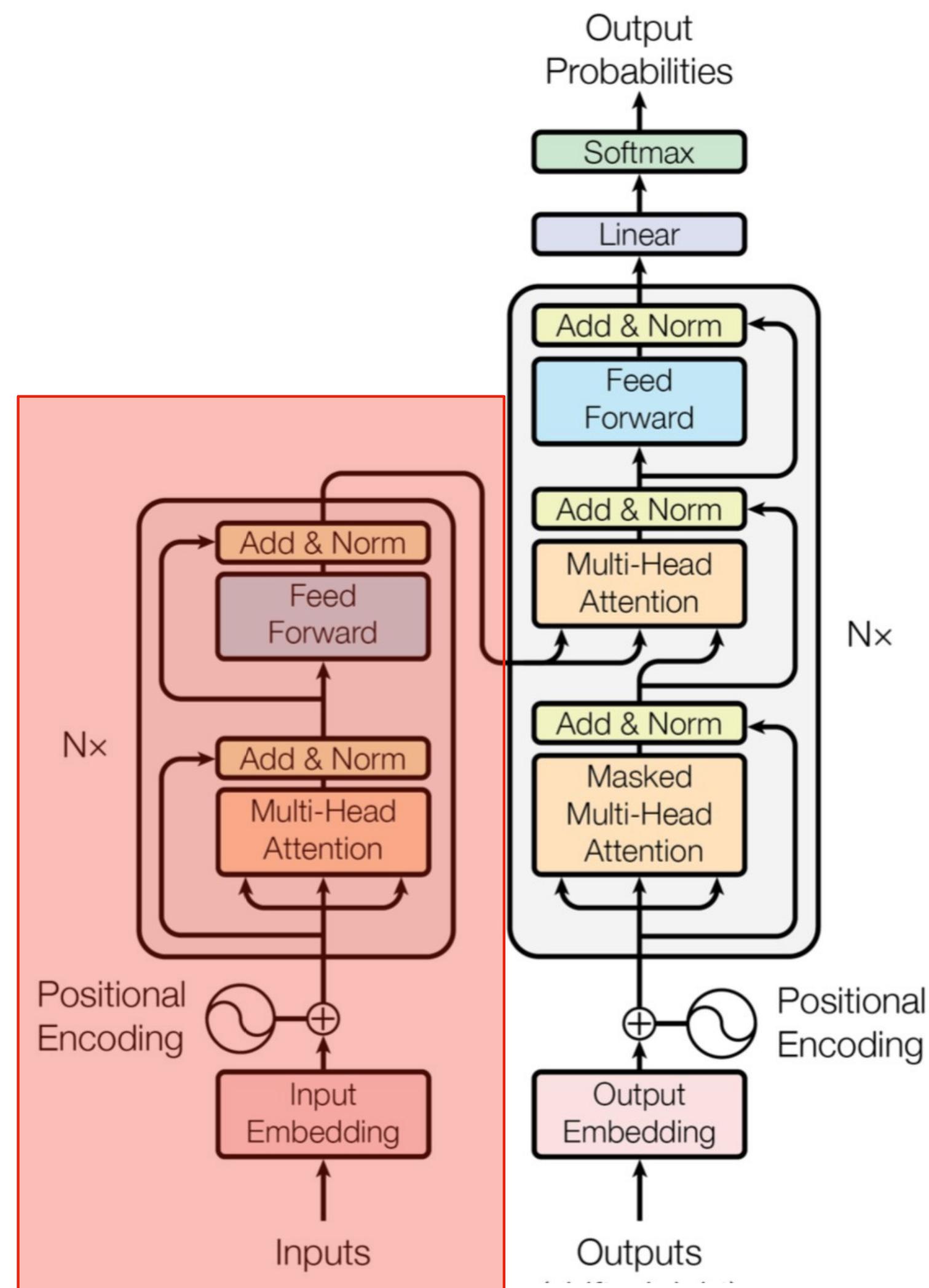
The Transformer Architecture



Attention is all you need (2017)

<https://arxiv.org/abs/1706.03762>

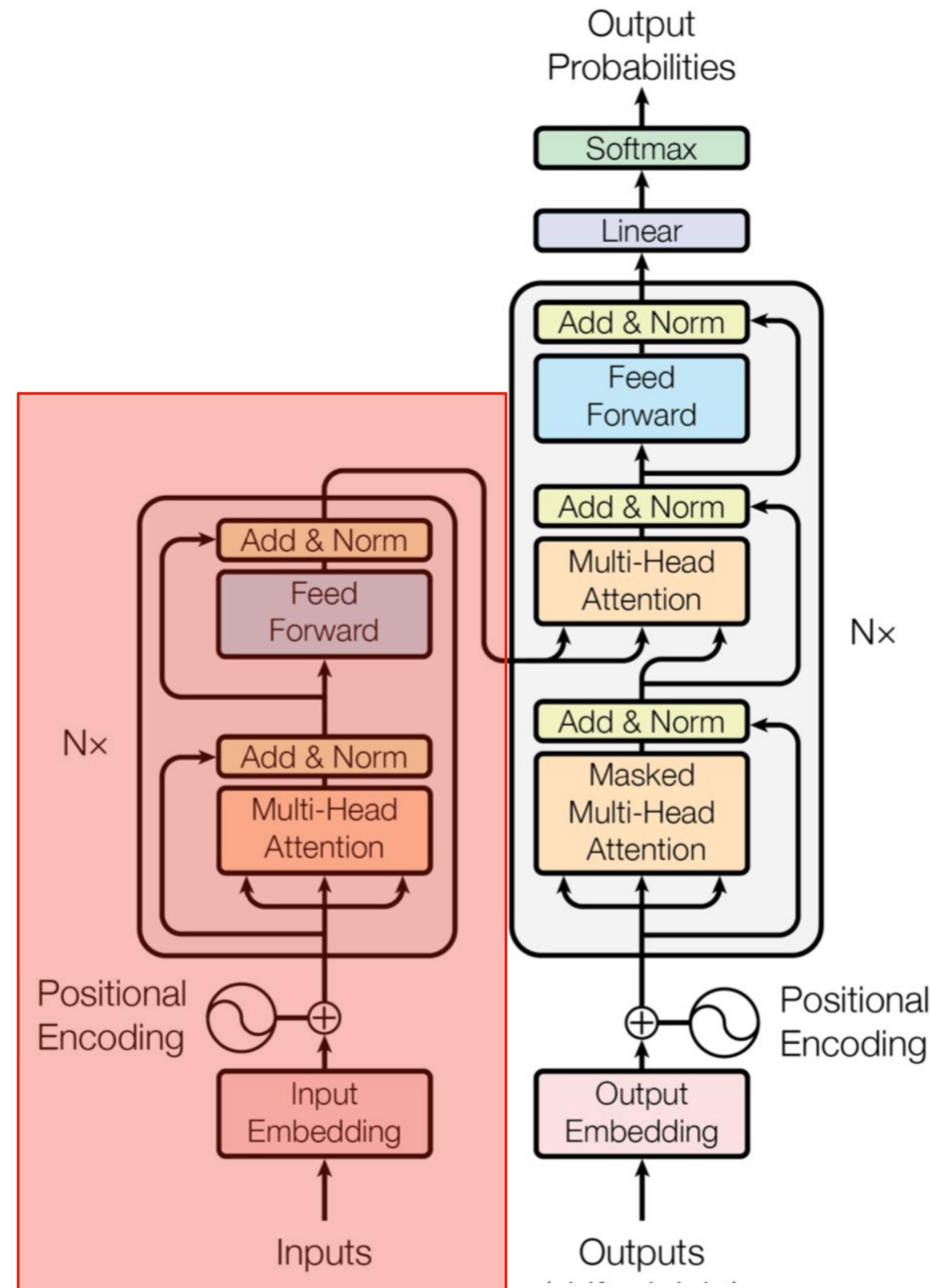
- Ground-breaking architecture that set SOTA on first translation and later all other NLP tasks
- For simplicity, can just look at one half of it



Transformer

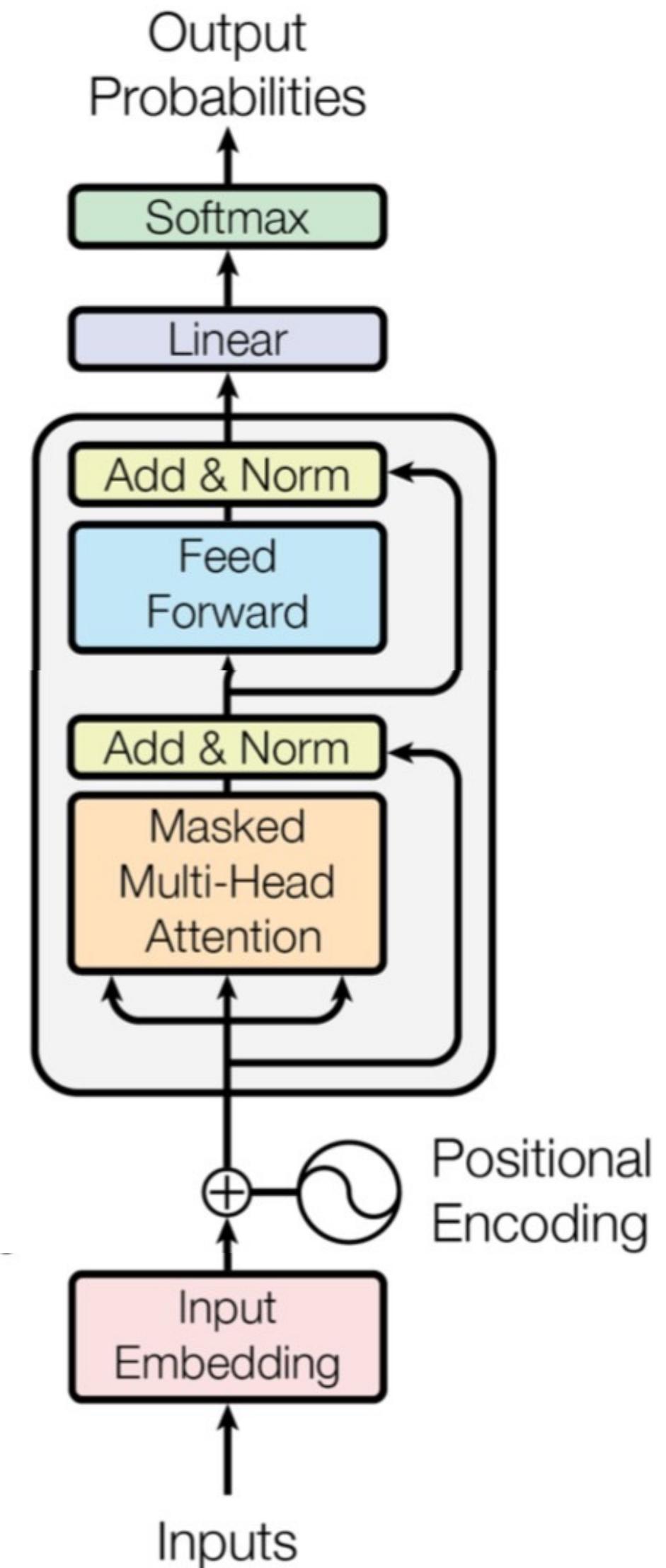
A breakthrough in AI that changes how we process data (i.e., text).

The core idea: Instead of reading word by word, transformers focus on the whole sentence and **important** parts of a sentence all at once.



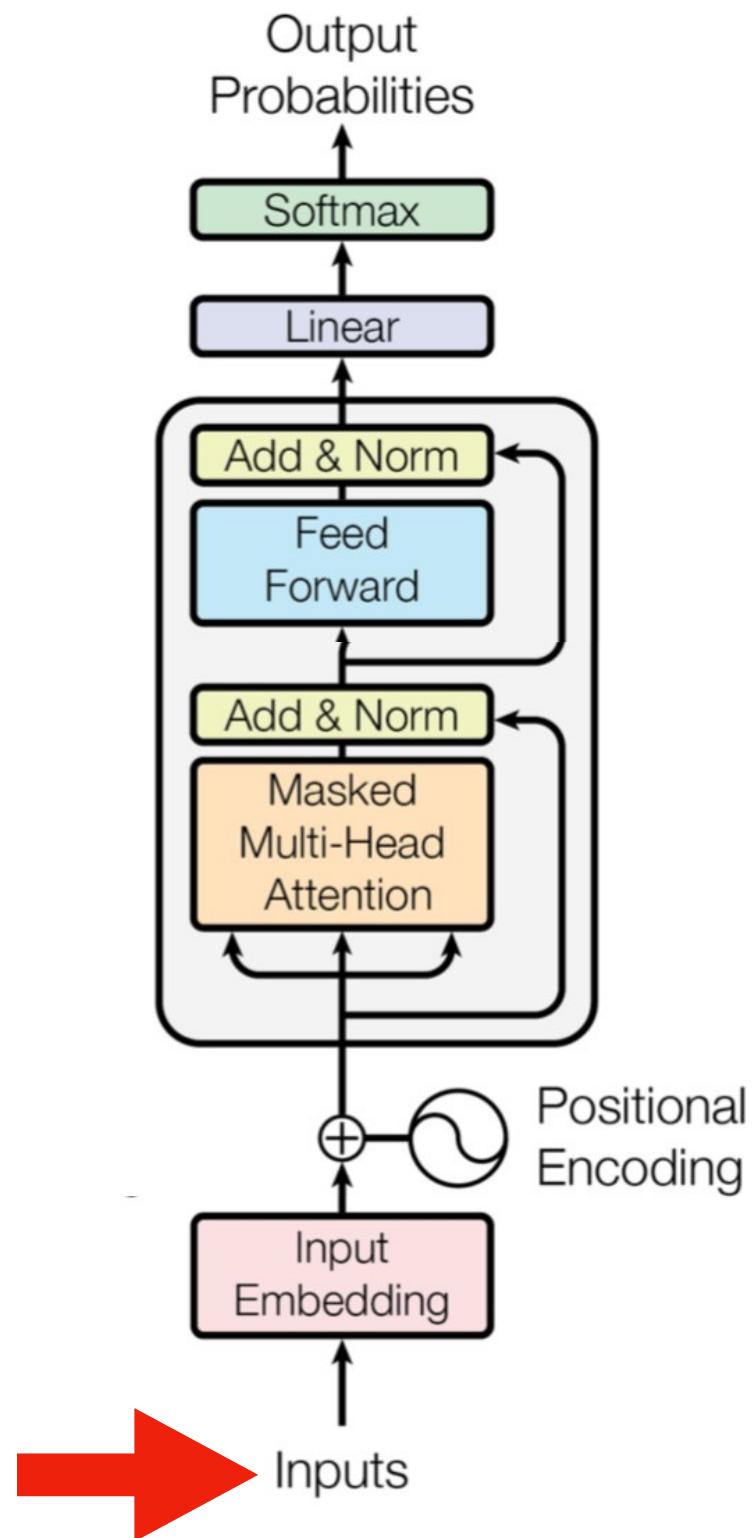
Transformer Decoder Overview

- Task is to complete text
 - "It's a blue" -> "sundress"
- Inputs: a sequence of N tokens
 - [It's, a, blue]
- Output:
 - Probability distribution over the next token
- Inference:
 - Sample the next token from the distribution, append it to inputs, run through the model again, sample, append, etc.



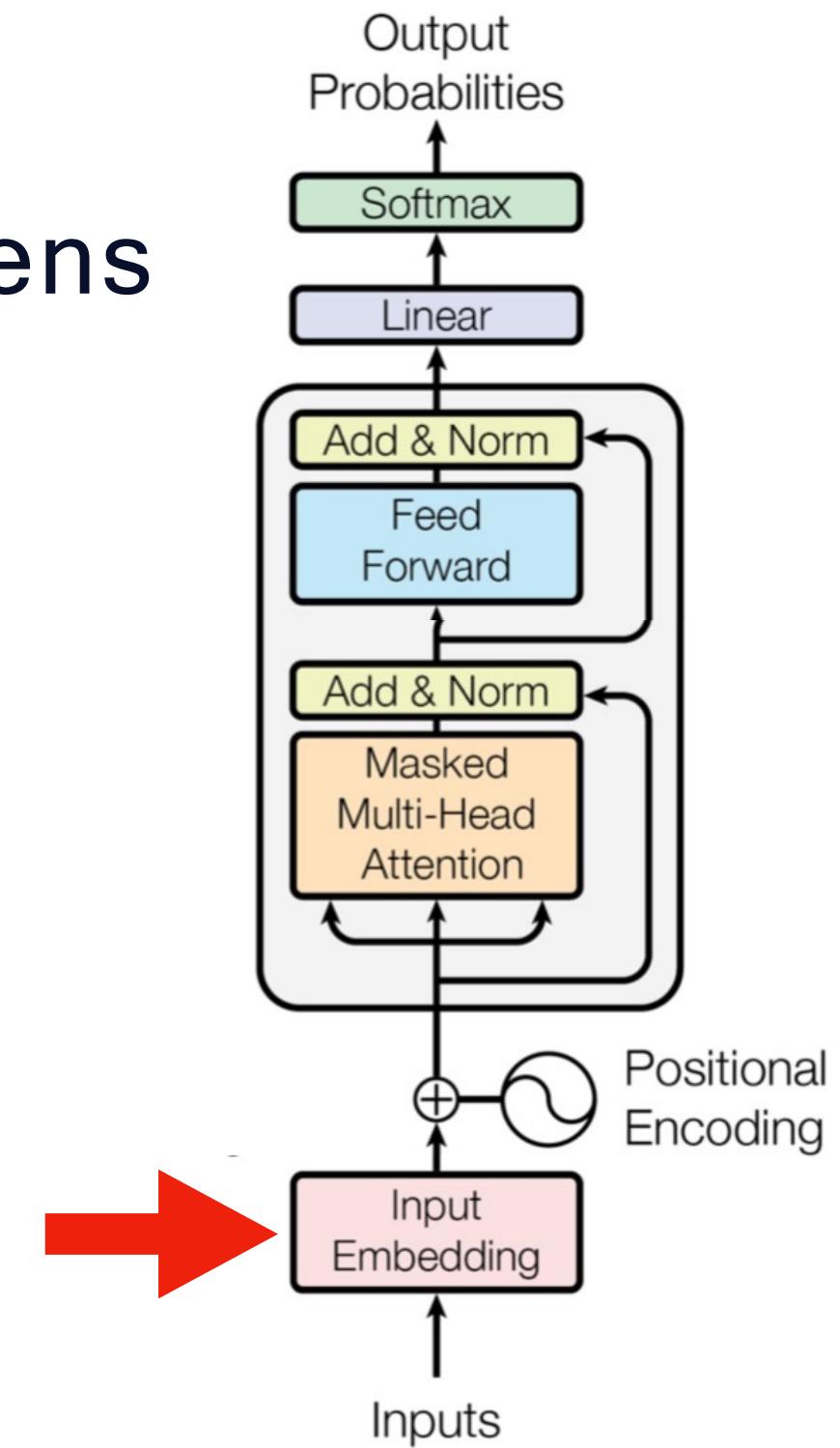
Vectorizing Inputs

- Inputs need to be vectors of numbers
 - Start with original text:
 - "It's a blue sundress."
- Turn into a sequence of tokens:
 - [<SOS>, It, 's, a, blue, sund, ress, ., <EOS>]
- Turn into vocabulary IDs:
 - [0, 1026, 338, 257, 4171, 37437, 601, 13, 1]
- Each ID can be represented by a one-hot vector
 - e.g. 13 -> [0, 0, 0, 1, 0, 0, 0, 0, ...] (all zeros, but one)



One-hot vectors

- One-hot vectors are poor representations of words or tokens
 - e.g. distance between "cat" and "kitten" is the same as between "cat" and "tractor"
- Cat: [1, 0, 0]
- Kitten: [0, 1, 0]
- Tractor: [0, 0, 1]



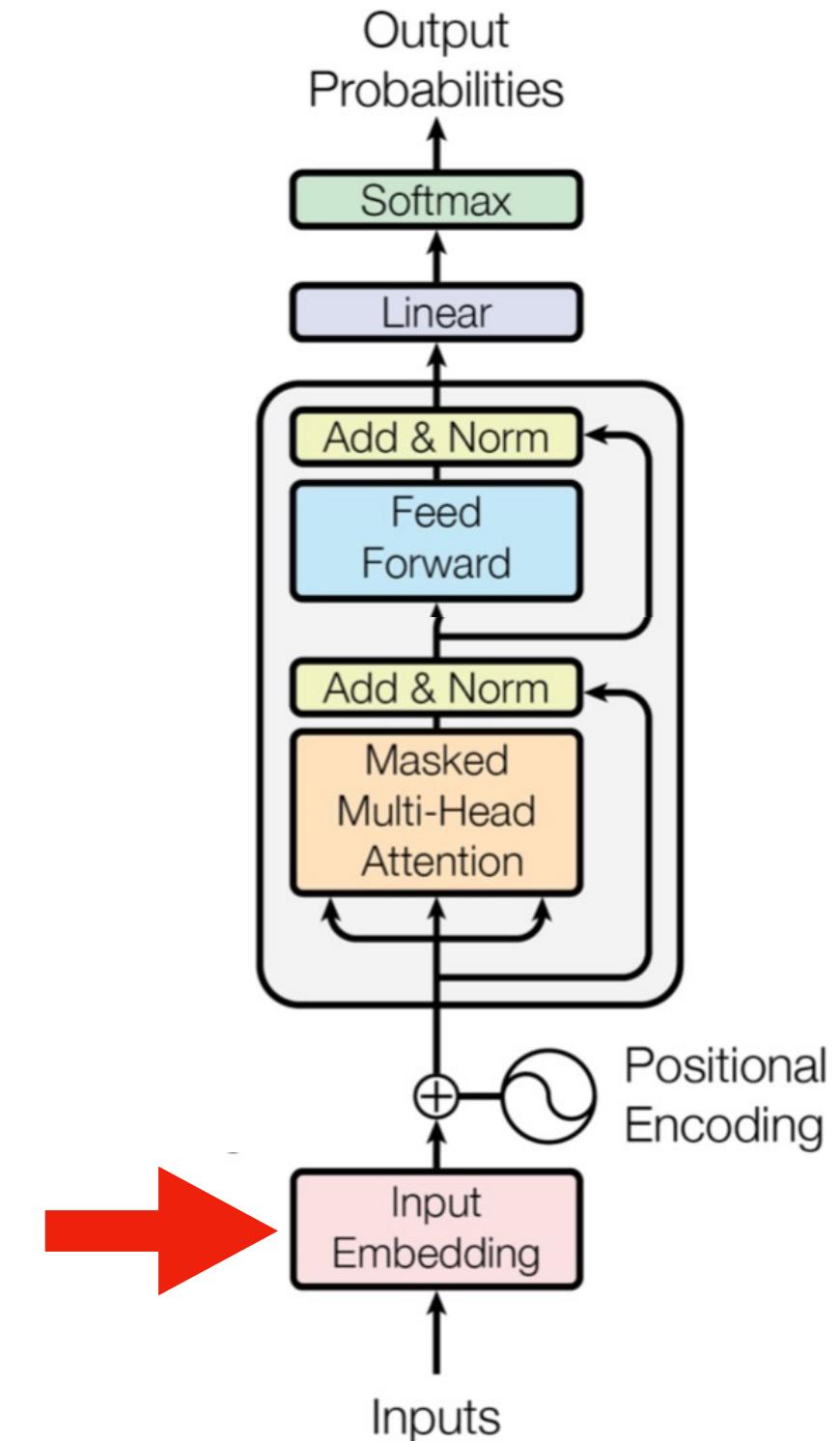
Input Embedding

- Solution: **learn** an embedding matrix!
 - An embedding maps each word to a high-dimensional vector of continuous numbers
 - learned through training a model on a specific task
 - e.g., Word2Vec: learns embeddings on large corpora of text

	1	0	0	...	0	0	0
aardvark	1	0	0	...	0	0	0
black	0	0	...	1	...	0	0
cat	0	0	...	1	...	0	0
duvet	0	0	...	1	...	0	0
zombie	0	0	0	...	0	0	1

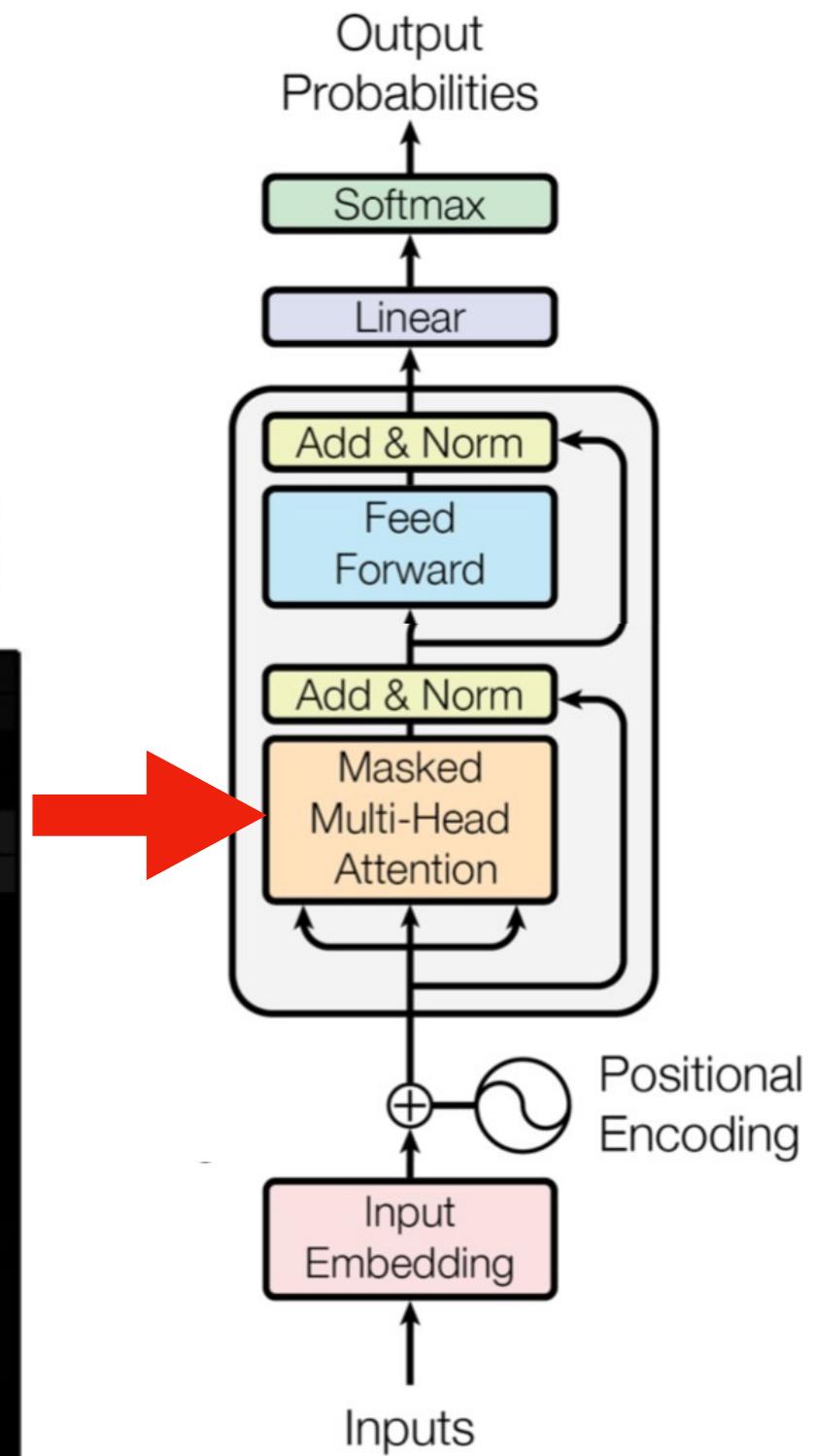
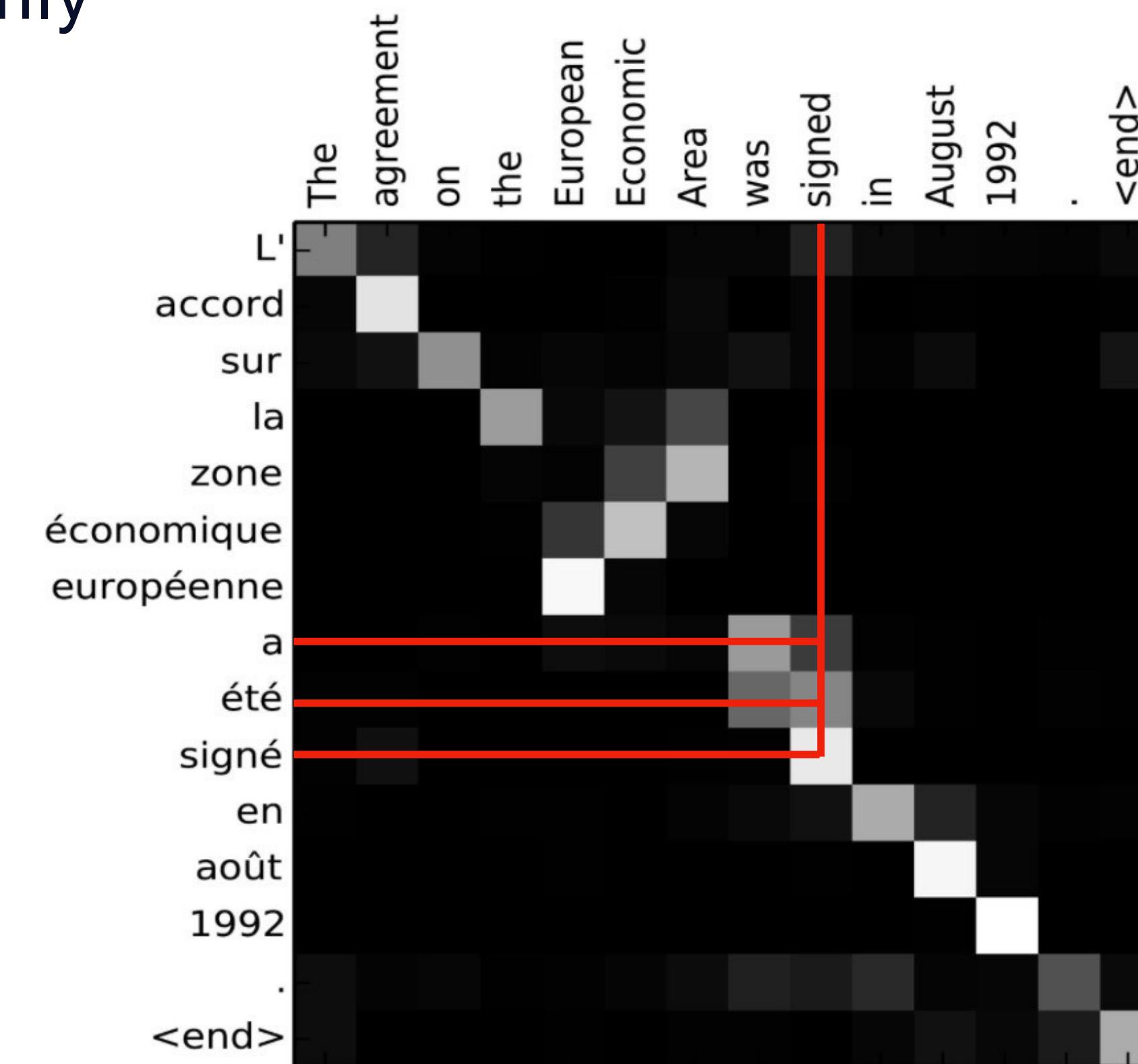
$V \times E$
*embedding
matrix*

aardvark	0.97	0.03	0.15	0.04
black	0.07	0.01	0.20	0.95
cat	0.98	0.98	0.45	0.35
duvet	0.01	0.84	0.12	0.02
zombie	0.74	0.05	0.98	0.93



Attention

- Key insight: for a given token in the output sequence, only one or a few tokens in the input sequence are most important
- Find those important relations
- Introduced in 2015 for translation tasks



Basic self-attention

- Input: sequence of vectors

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t$$

- Output: sequence of vectors, each one a weighted sum of the input sequence

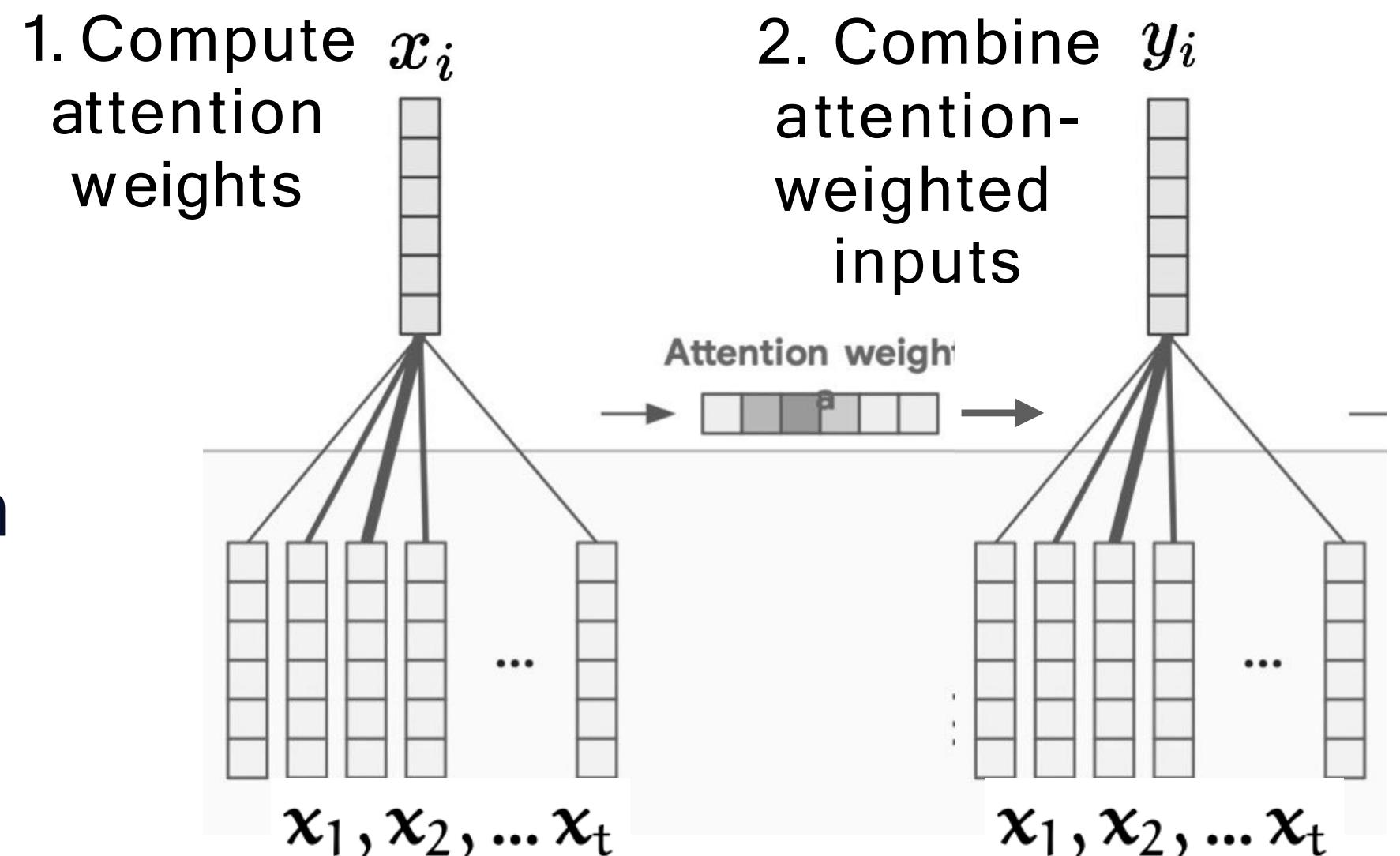
$$\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_t$$

$$\mathbf{y}_i = \sum_j w_{ij} \mathbf{x}_j$$

- weight is just dot product between input vectors

- (made to sum to 1)

$$w'_{ij} = \mathbf{x}_i^\top \mathbf{x}_j$$
$$w_{ij} = \frac{\exp w'_{ij}}{\sum_j \exp w'_{ij}}$$



Basic self-attention

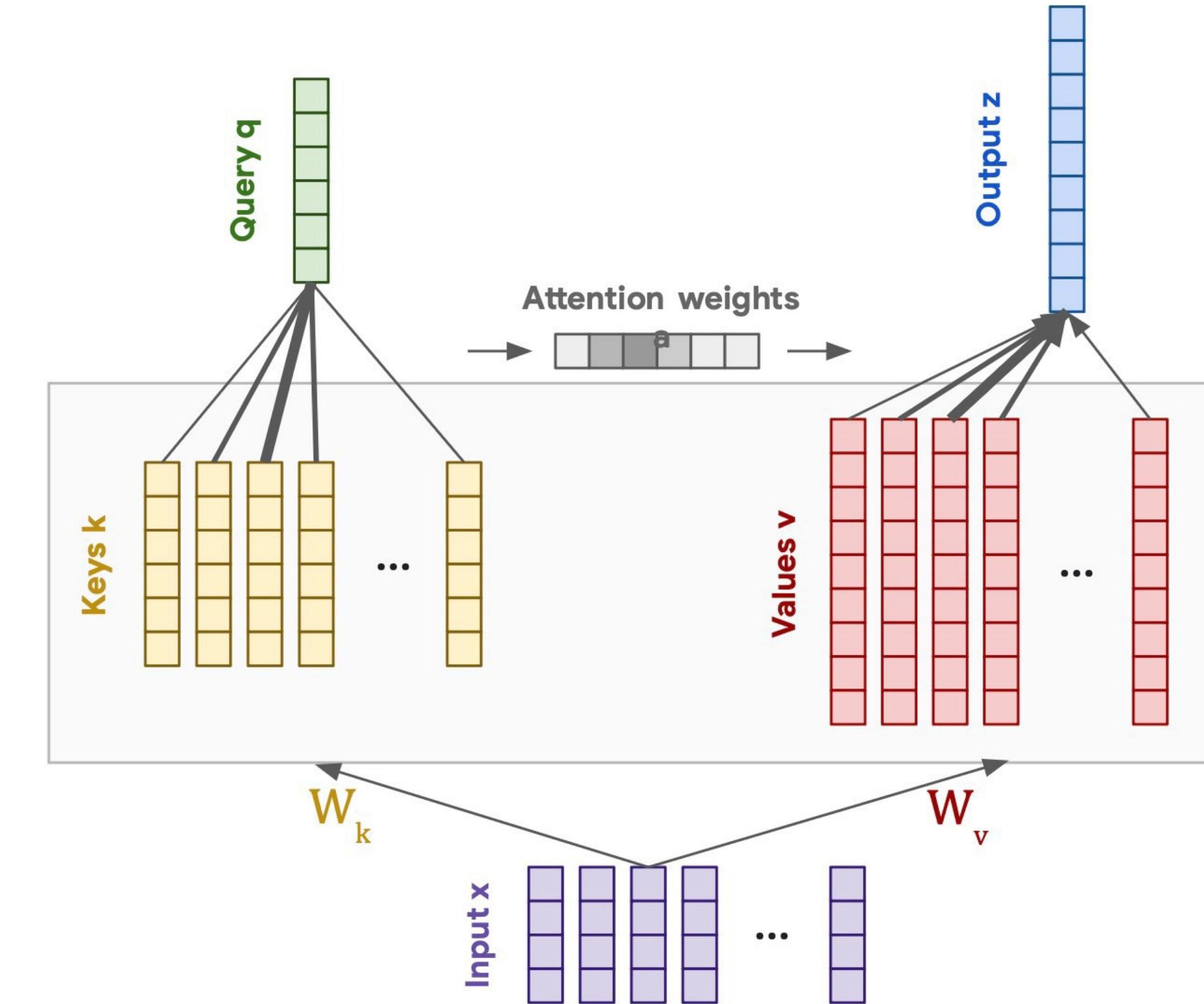
- Problem: there's no learning involved!
- Solution: project inputs into query, key, value roles
- Learning these matrices = learning attention

$$q_i = \mathbf{W}_q x_i \quad k_i = \mathbf{W}_k x_i \quad v_i = \mathbf{W}_v x_i$$

$$w'_{ij} = q_i^T k_j$$

$$w_{ij} = \text{softmax}(w'_{ij})$$

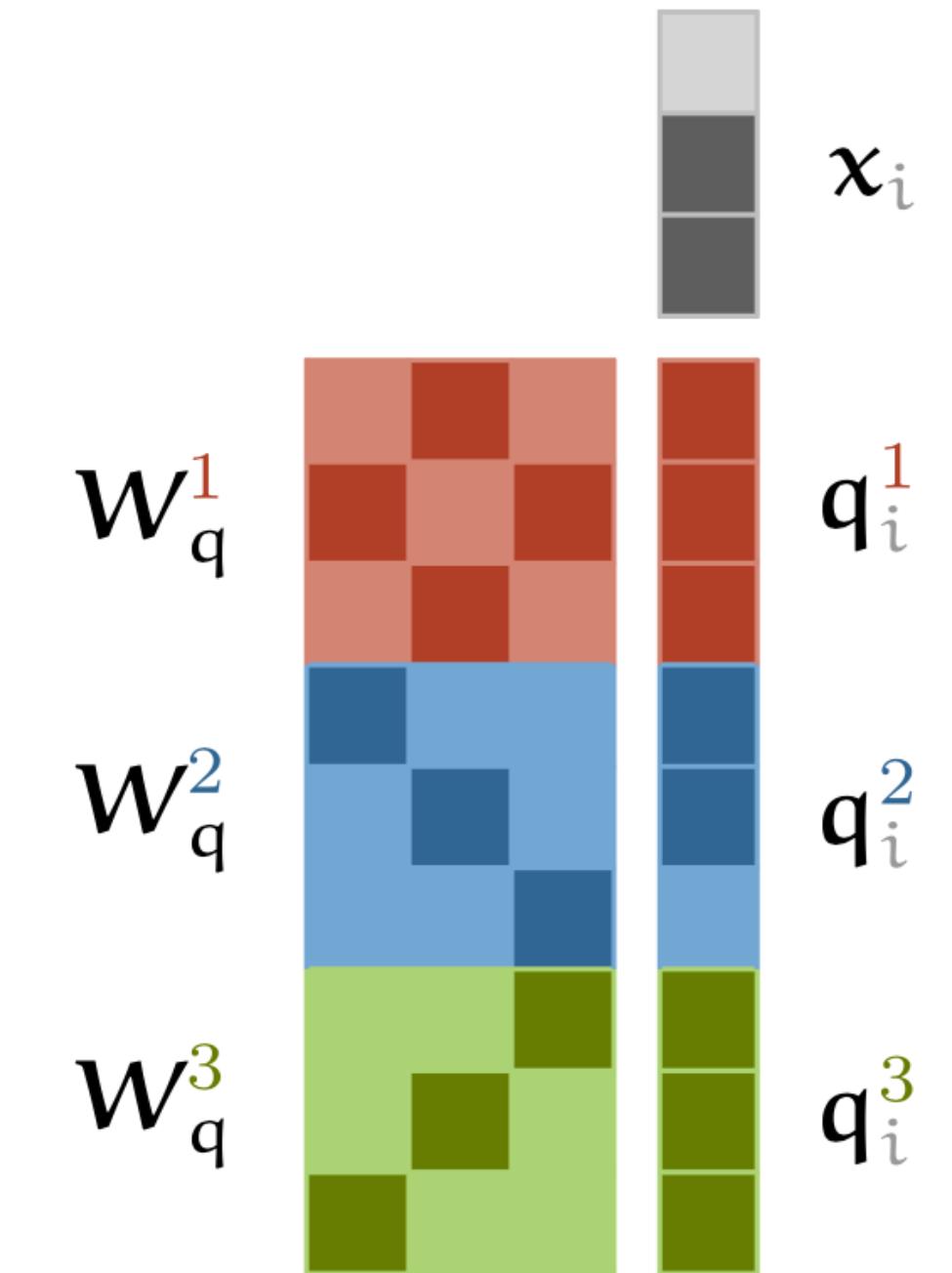
$$y_i = \sum_j w_{ij} v_j .$$



<http://lucasbeyer.be/transformer>

Multi-head attention

- Different parts of the sentence could have different relations
- We can allow different ways of transforming into queries, keys, and values to be learned
- Simply means learning different sets of W_q , W_k , and W_v matrices simultaneously.

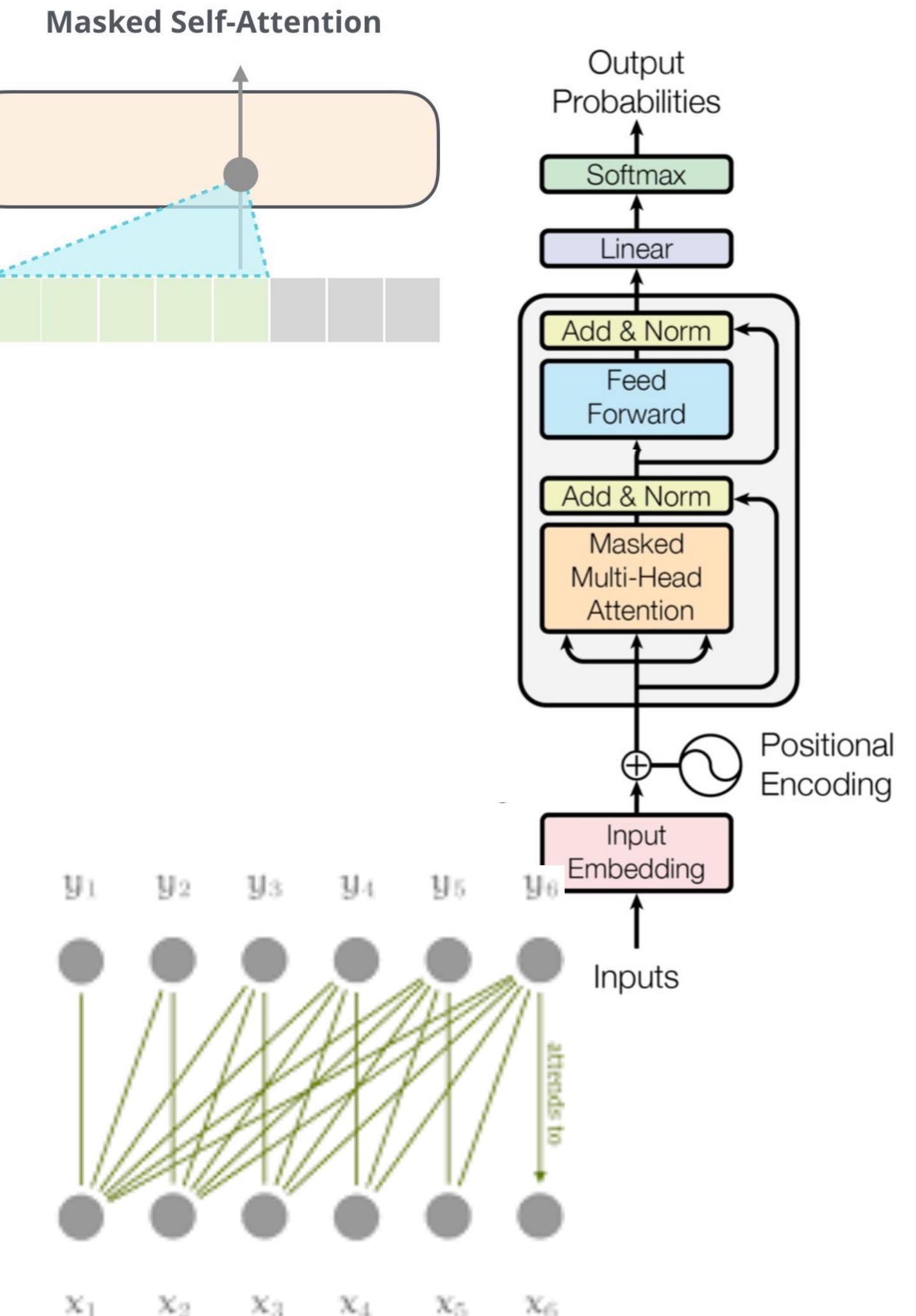
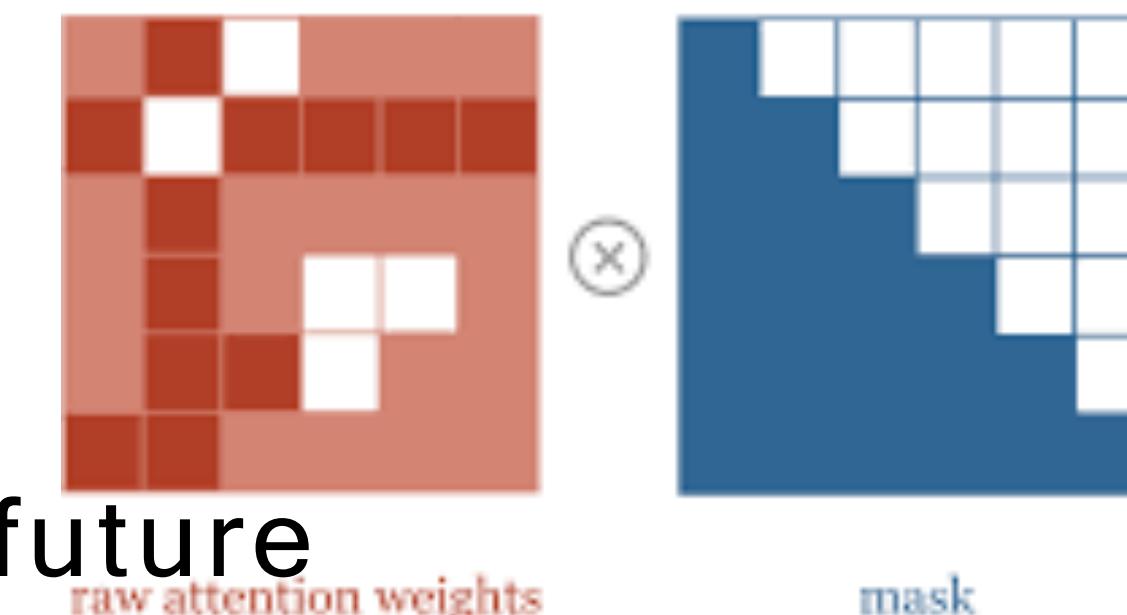


3-headed attention

Masking attention

In training:

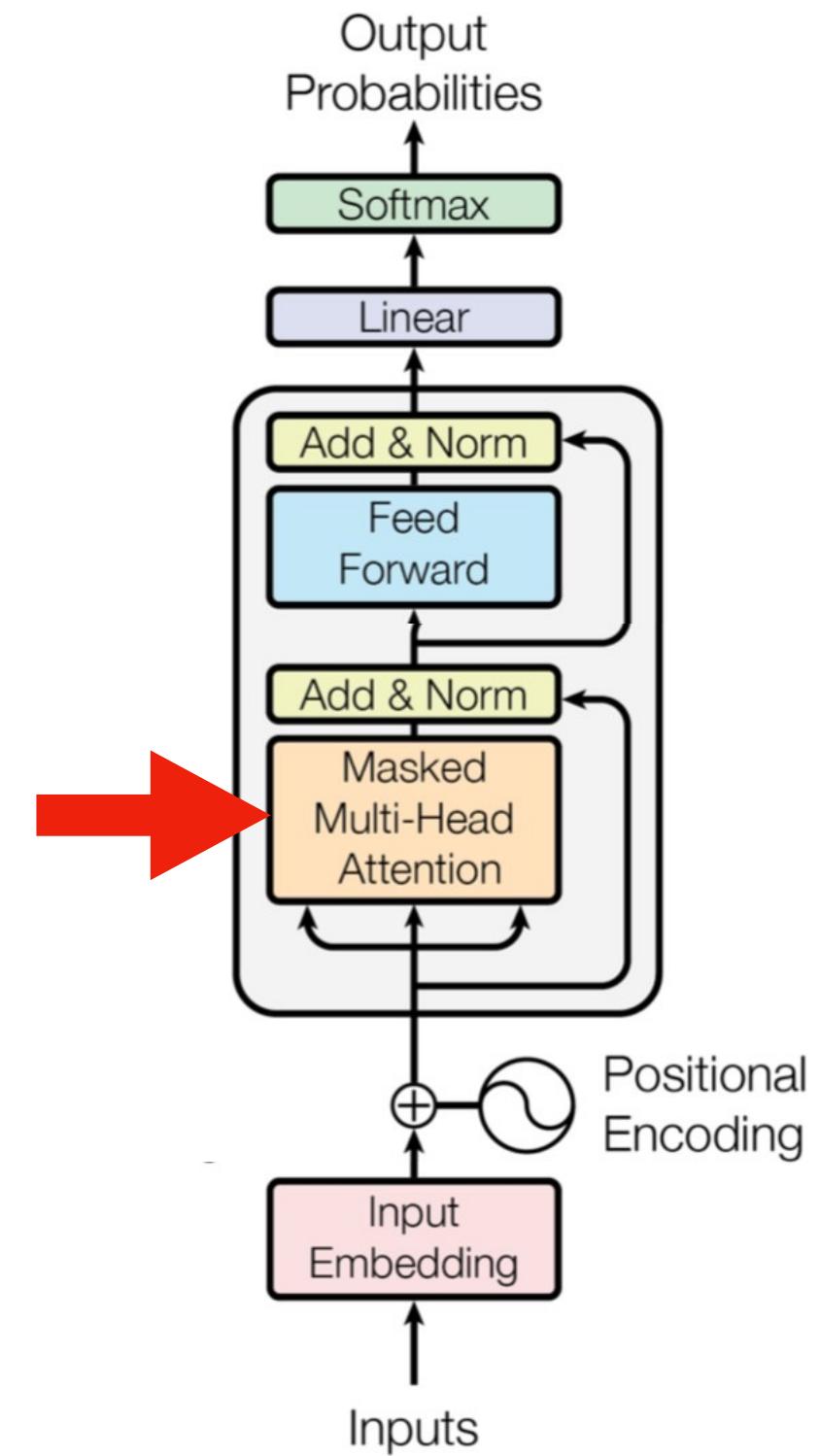
- Inputs: a sequence of N tokens
 - [It's, a, blue, <BLANK>, <BLANK>, ...]
- Ground-truth Outputs: a sequence of N tokens
 - [a, blue, sundress, <BLANK>, <BLANK>, ...]
- Actual Outputs: Instead of words, vectors of probabilities over the vocabulary
 - **Crucially: all output probabilities are computed at the same time!**



Note how you shouldn't see future tokens when predicting

Masked Multi-Head Attention

- Conceptual view:
 - token comes in
 - gets "augmented" with previously-seen tokens that seem relevant (masked self-attention)
 - this happens in several different ways simultaneously (multiple heads)
- NOTE: there's no notion of "position" so far!

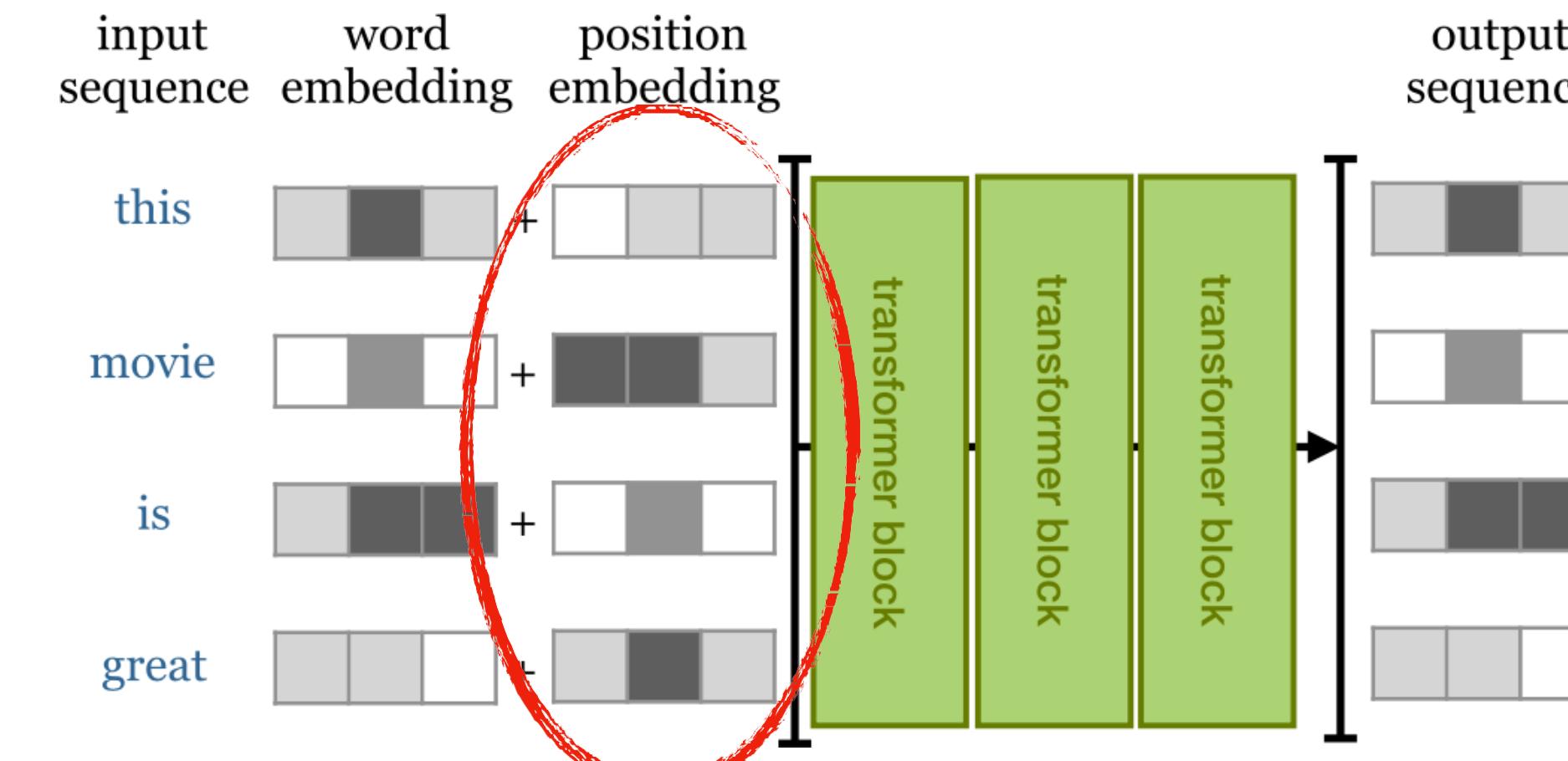
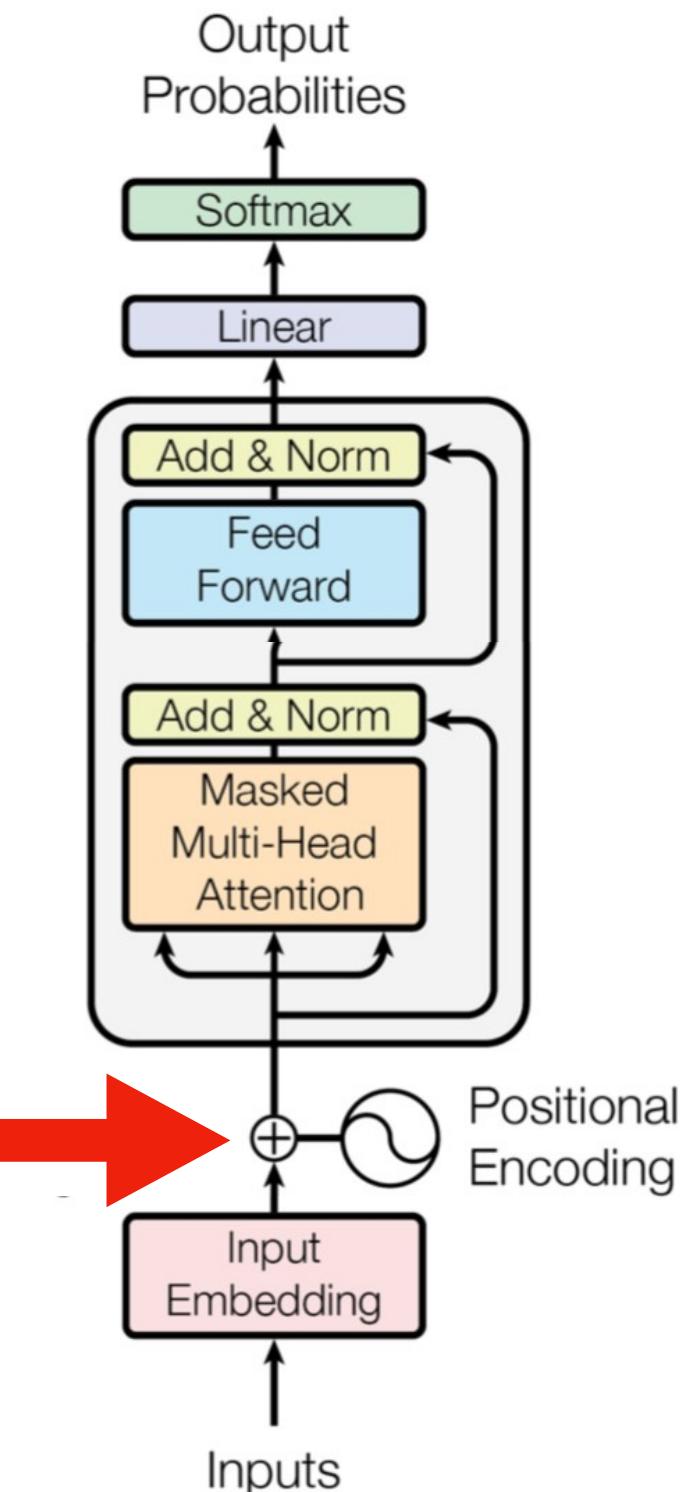


$$\mathbf{q}_i = \mathbf{W}_q \mathbf{x}_i \quad \mathbf{k}_i = \mathbf{W}_k \mathbf{x}_i \quad \mathbf{v}_i = \mathbf{W}_v \mathbf{x}_i$$

Positional Encoding

- Attention is totally position-invariant!
 - e.g. [this, movie, is, great] is the same as [movie, this, great, is]
- So, let's add position-encoding vectors to embedding vectors
 - It really is that simple

$$\begin{aligned} w'_{ij} &= \mathbf{q}_i^\top \mathbf{k}_j \\ w_{ij} &= \text{softmax}(w'_{ij}) \\ \mathbf{y}_i &= \sum_j w_{ij} \mathbf{v}_j . \end{aligned}$$

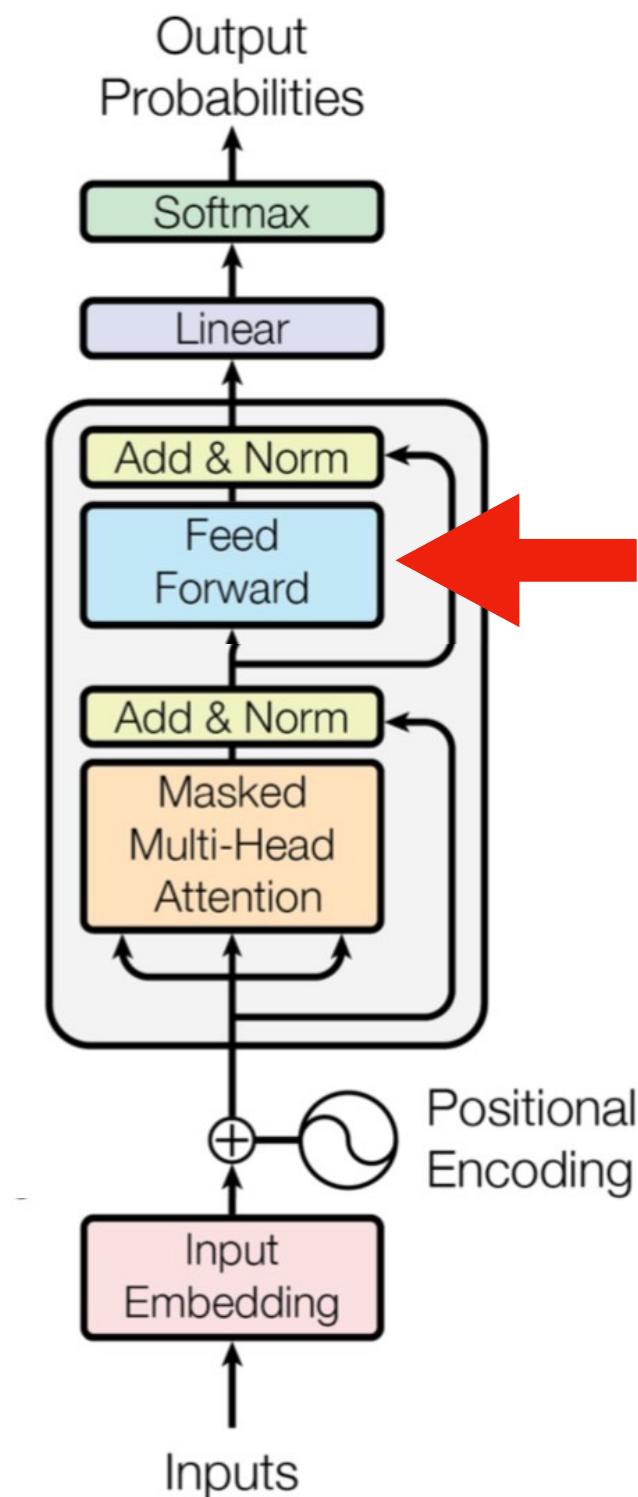
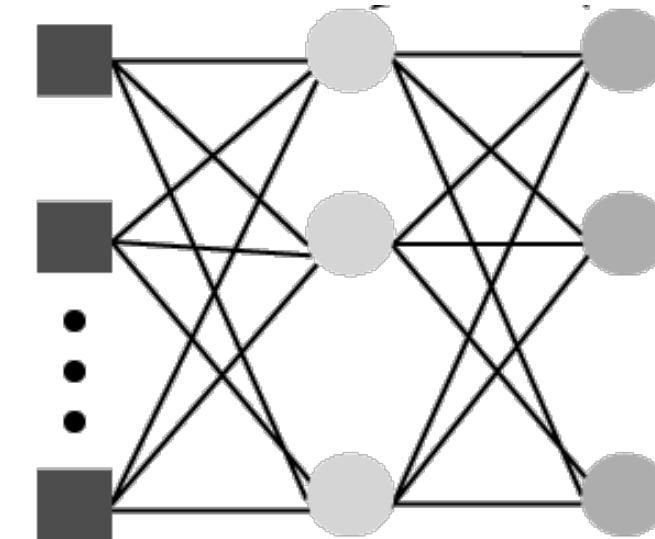


Feed Forward Layer

- Standard Multi-Layer Perceptron with one hidden layer

- Conceptual view:

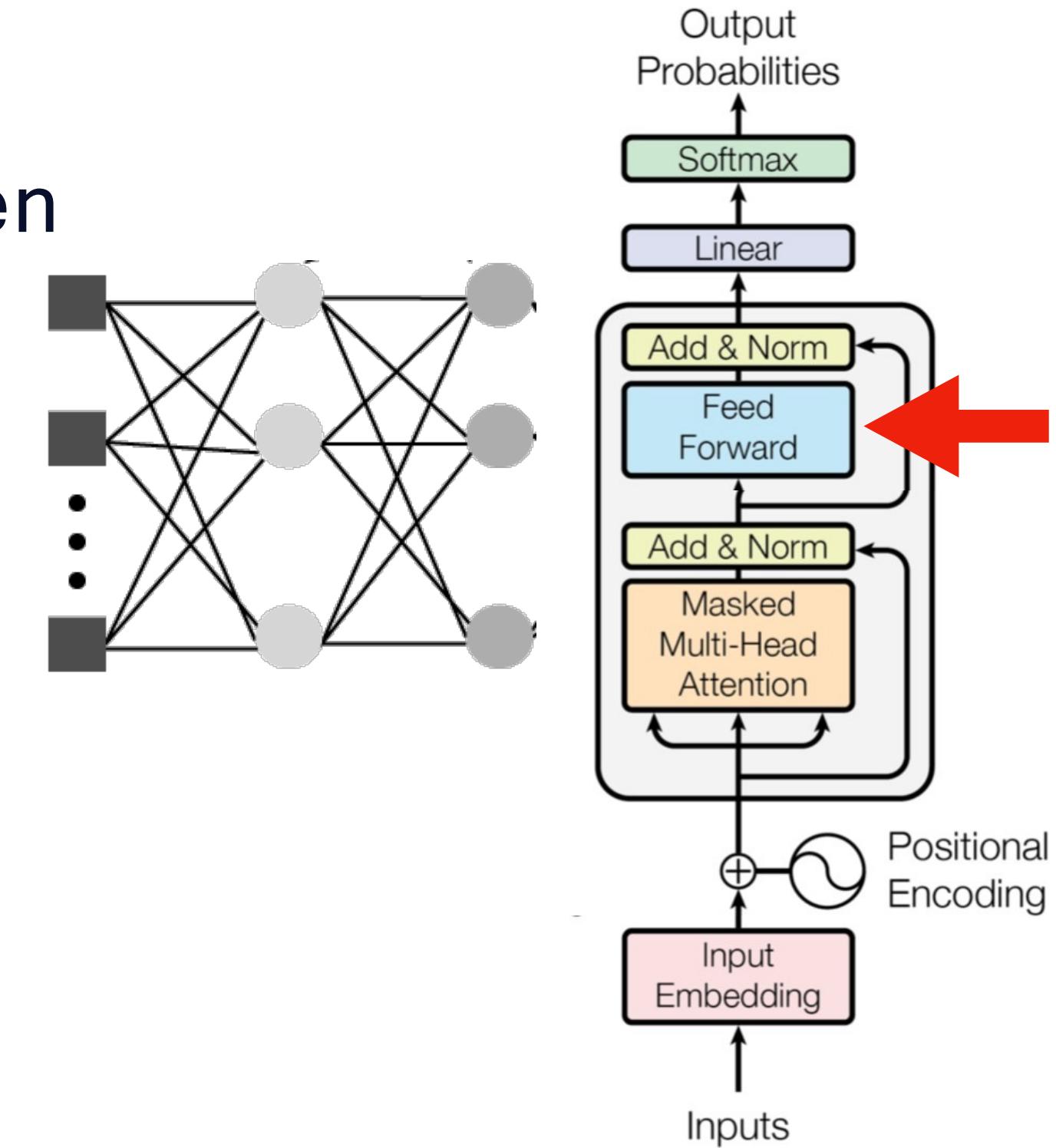
- token (augmented with other relevant tokens that it has seen) comes in...
- ...and "upgrades" its representation



Feed Forward Layer

- Standard Multi-Layer Perceptron with one hidden layer

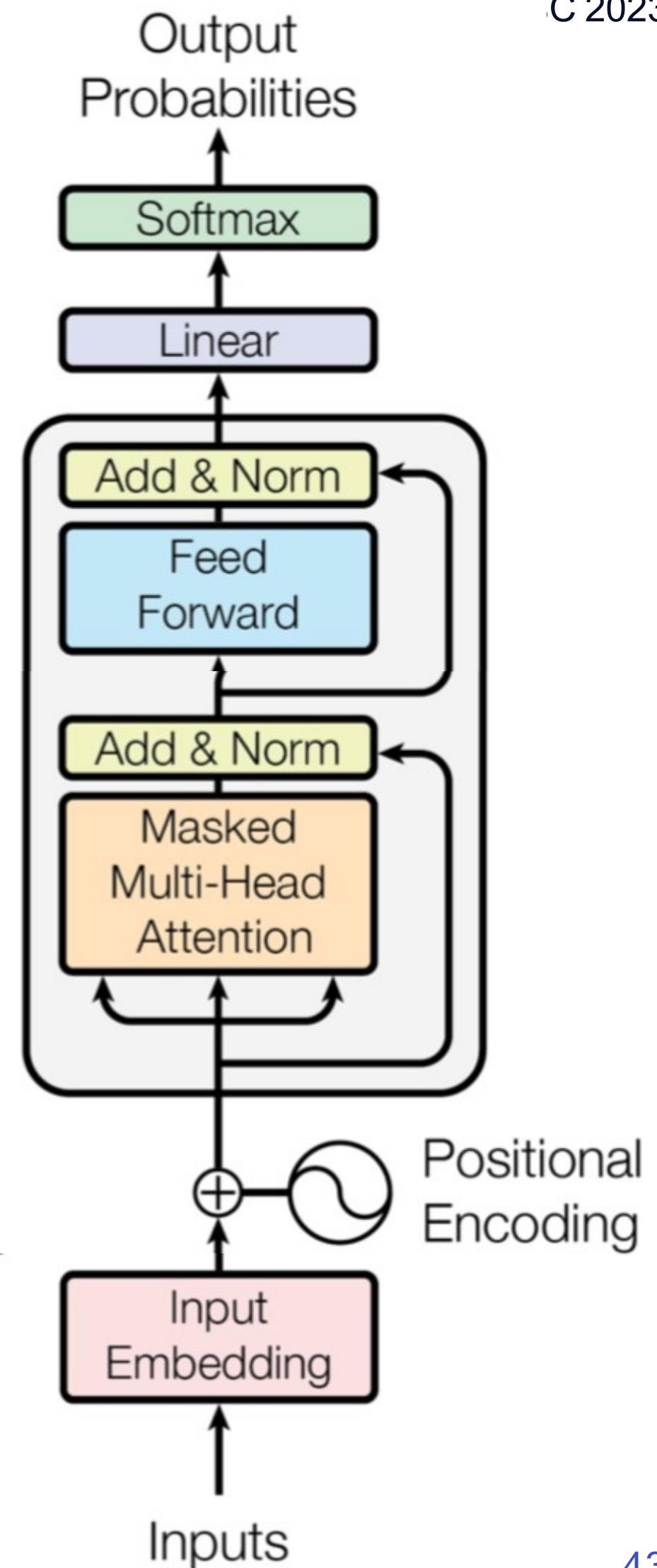
After the attention mechanism aggregates and weighs information across the input sequence, the feed-forward network applies a series of linear transformations to this aggregated information. This process allows the model to learn complex relationships and patterns beyond what the attention mechanism can do on its own.



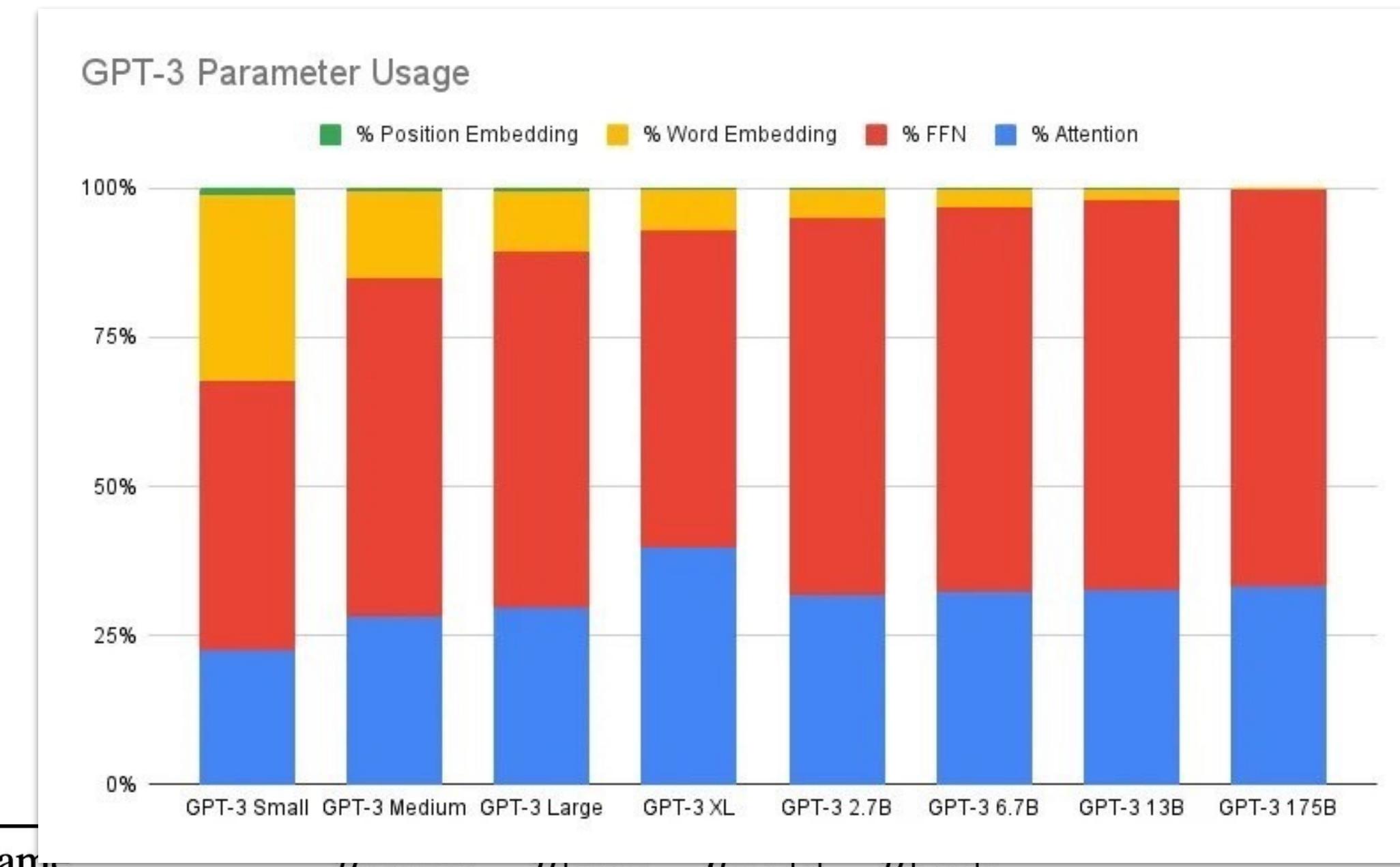
Transformer Architecture

- The main Transformer Layer is stacked many times
- The overall hyperparameters are:
 - Number of layers
 - Embedding dimension
 - Number of attention heads
- The largest models are ~70% feed-forward weights

<https://aizi.substack.com/p/how-does-gpt-3-spend-its-175b-parameters>



Transformer Architecture



Model Name	#params	#layers	#model	#heads
GPT-3 Small	125M	12	768	12
GPT-3 Medium	350M	24	1024	16
GPT-3 Large	760M	24	1536	16
GPT-3 XL	1.3B	24	2048	24
GPT-3 2.7B	2.7B	32	2560	32
GPT-3 6.7B	6.7B	32	4096	32
GPT-3 13B	13.0B	40	5140	40
GPT-3 175B or “GPT-3”	175.0B	96	12288	96



Why does this work so well?



Andrej Karpathy 

@karpathy

...

The Transformer is a magnificent neural network architecture because it is a general-purpose differentiable computer. It is simultaneously:

- 1) expressive (in the forward pass)
- 2) optimizable (via backpropagation+gradient descent)
- 3) efficient (high parallelism compute graph)

11:54 AM · Oct 19, 2022

490 Retweets

39 Quotes

3,670 Likes

1,023 Bookmarks

We mostly don't understand it, though

In-context Learning and Induction Heads

AUTHORS

Catherine Olsson*, Nelson Elhage*, Neel Nanda*, Nicholas Joseph†, Nova DasSarma†, Tom Henighan†, Ben Mann†, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, Chris Olah*

AFFILIATION

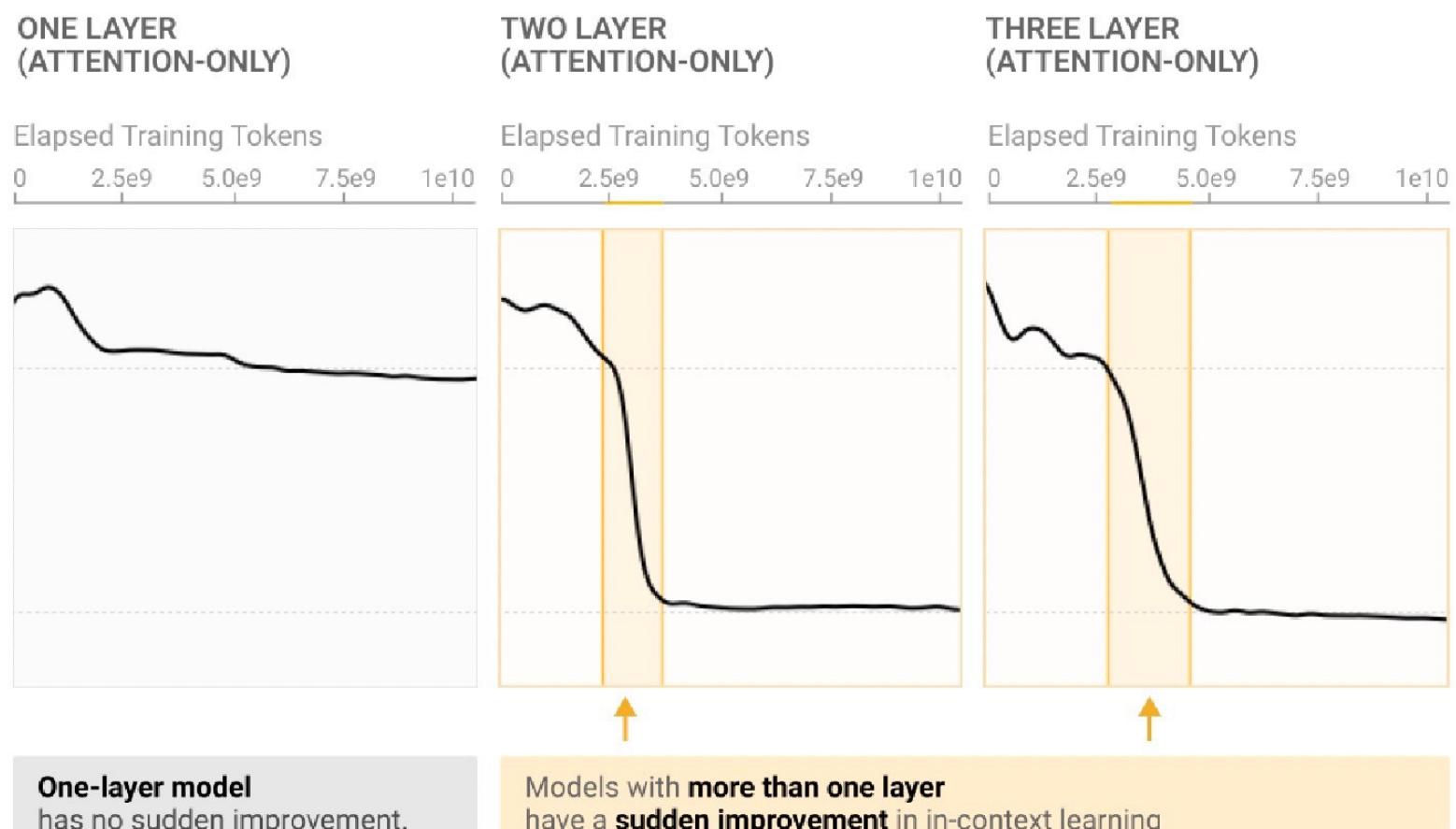
Anthropic

PUBLISHED

Mar 8, 2022

* Core Research Contributor; † Core Infrastructure Contributor; * Correspondence to colah@anthropic.com;
Author contributions statement below.

MODELS WITH MORE THAN ONE LAYER HAVE AN ABRUPT IMPROVEMENT IN IN-CONTEXT LEARNING

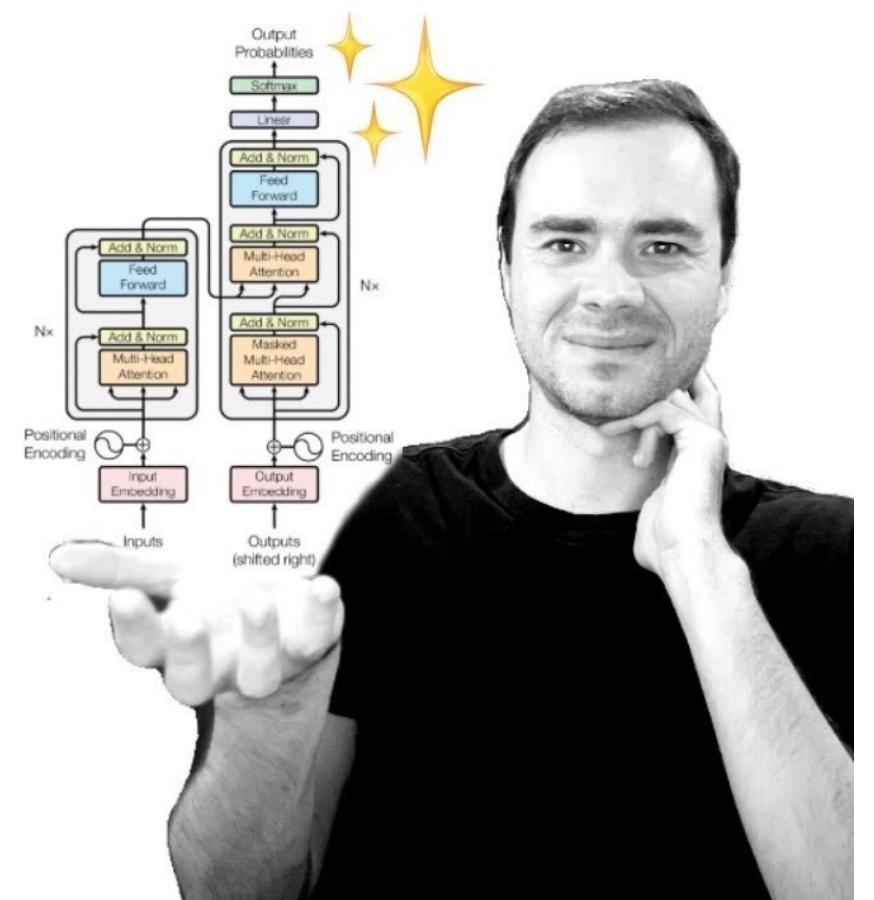


We highlight
“phase change”
period of training in
make visual compar
between plots easier.
highlighted region is
selected for each mo
based on the derivat
in-context learning.

Should you be able to code a Transformer?

- Definitely not necessary!
- BUT: it's not difficult, it is fun, and is probably worth doing
- Andrej Karpathy's GPT-2 implementation is <400 lines of code, including Attention

**LET'S BUILD GPT.
FROM SCRATCH.
IN CODE.
SPELLED OUT.**



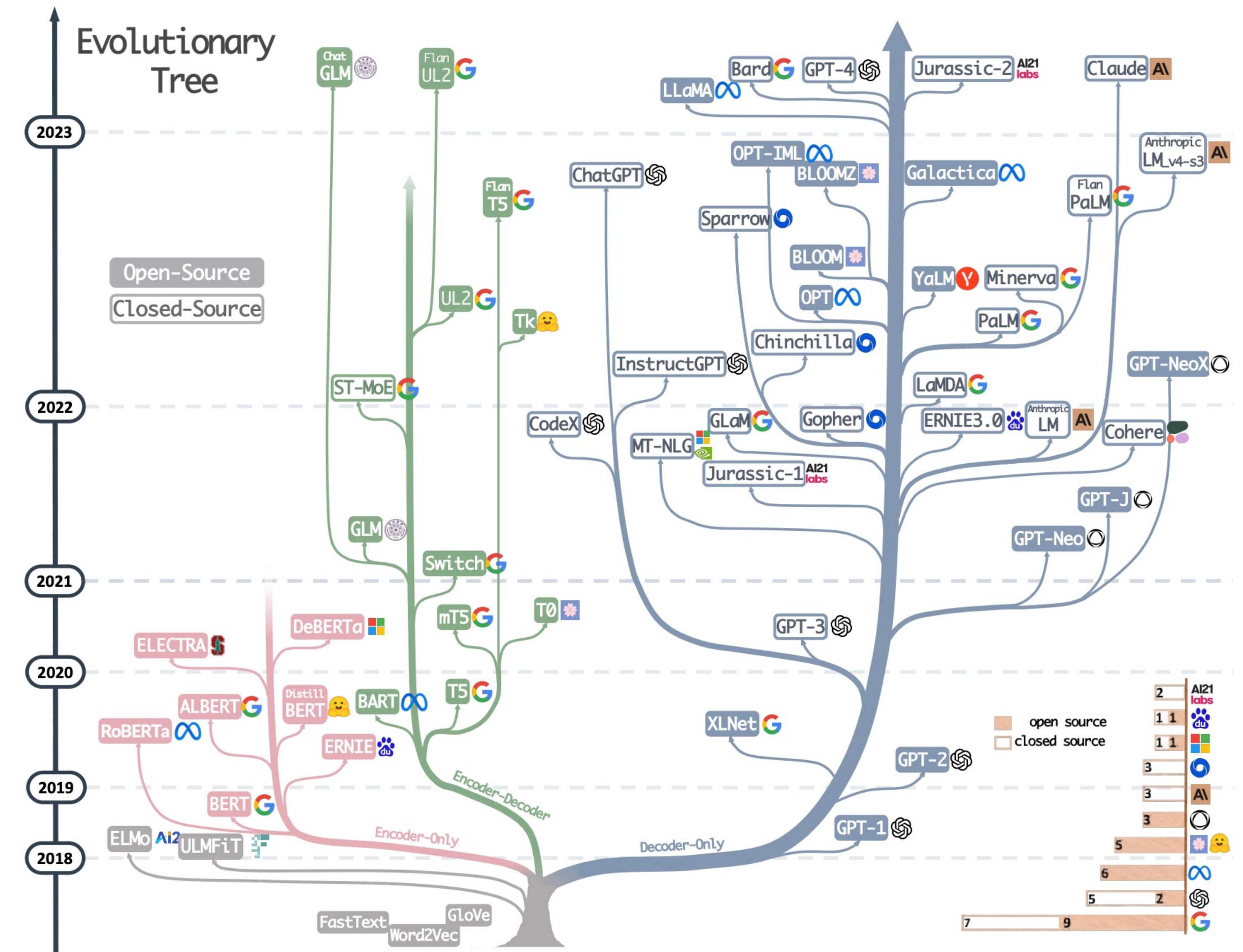
<https://www.youtube.com/watch?v=kCc8FmEb1nY>

Resources

- Lucas Beyer's [Lecture on Transformers](#)
- Peter Bloem's "[Transformers from Scratch](#)"
- Nelson Elhage's "[Transformers for Software Engineers](#)" for a different view
- Andrej Karpathy's entire [Neural Networks: Zero to Hero](#) video series
- Lillian Weng's "[The Transformer Family v2](#)" megapost

Questions?

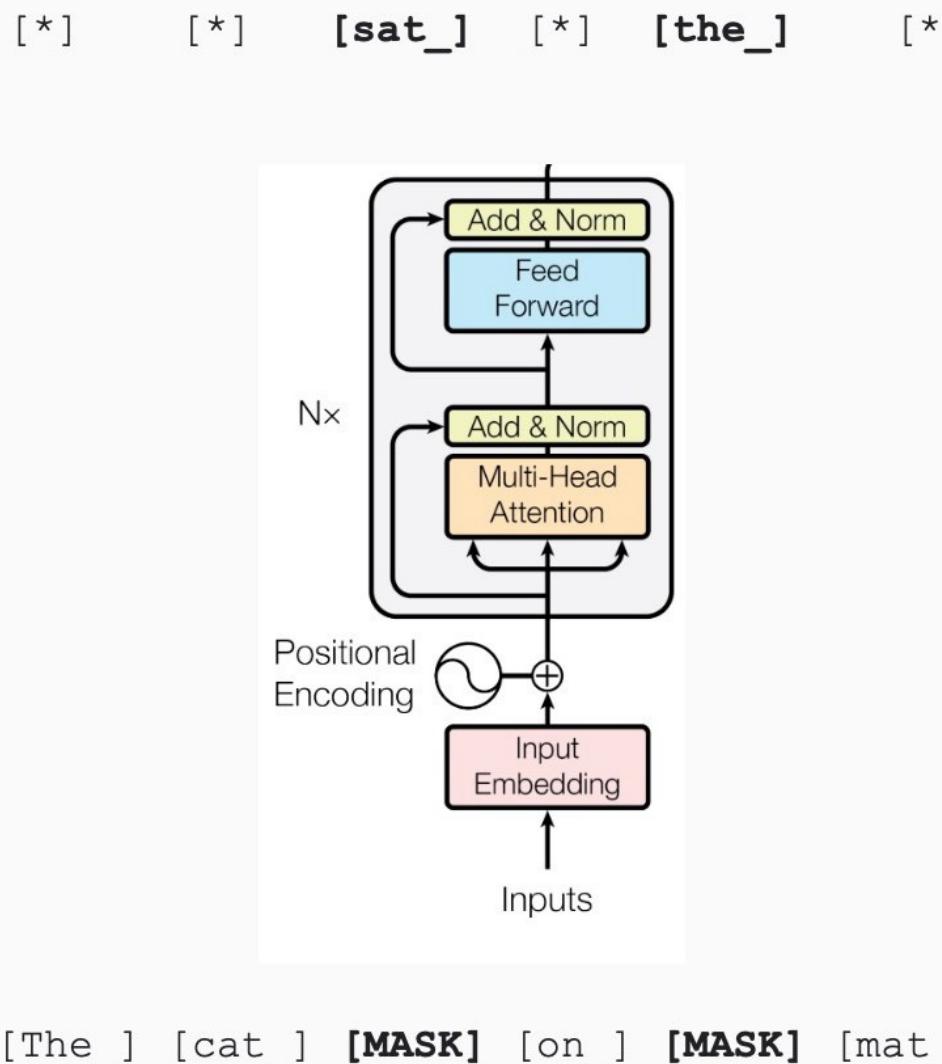
Notable LLMs





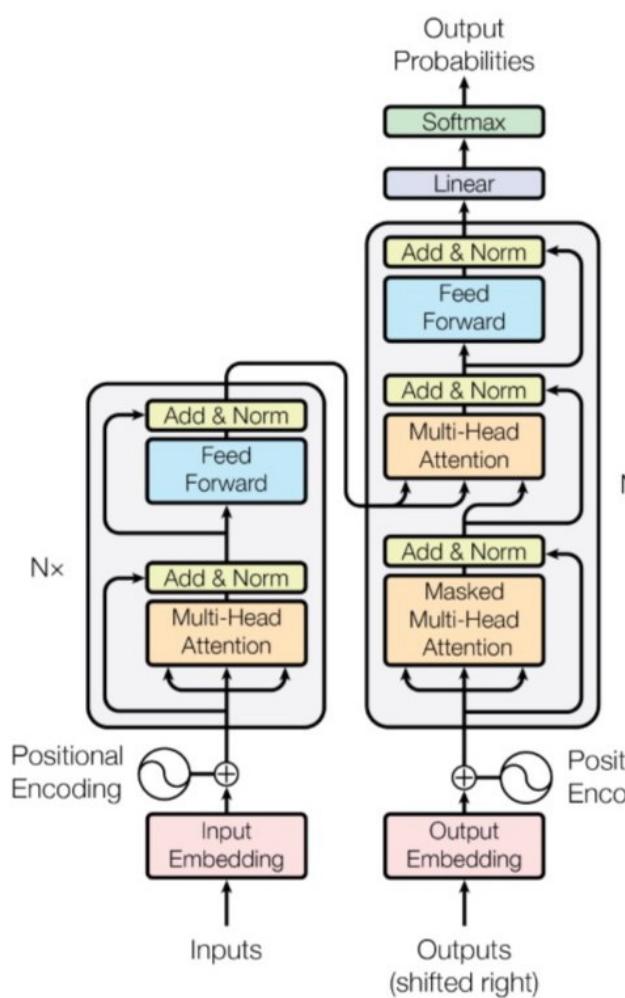
Three Easy Pieces

BERT



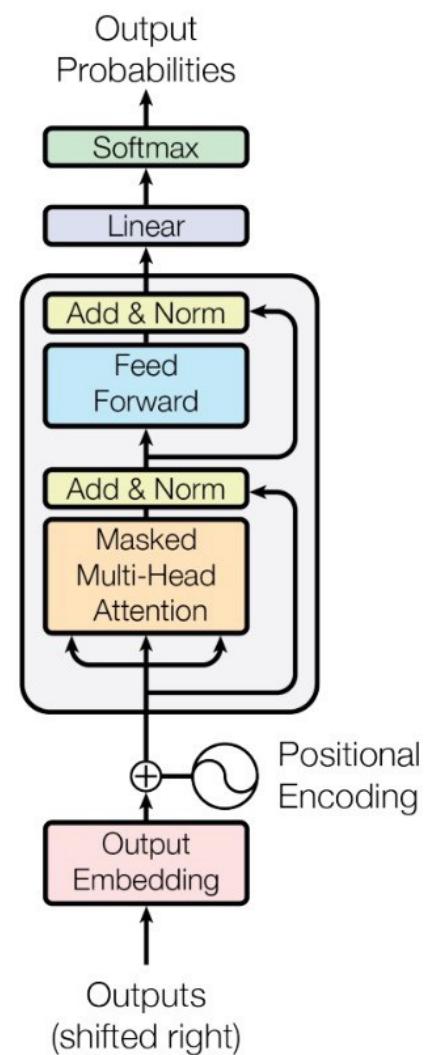
T5

Das ist gut.
A storm in Attala caused 6 victims.
This is not toxic.



GPT

[sat_]



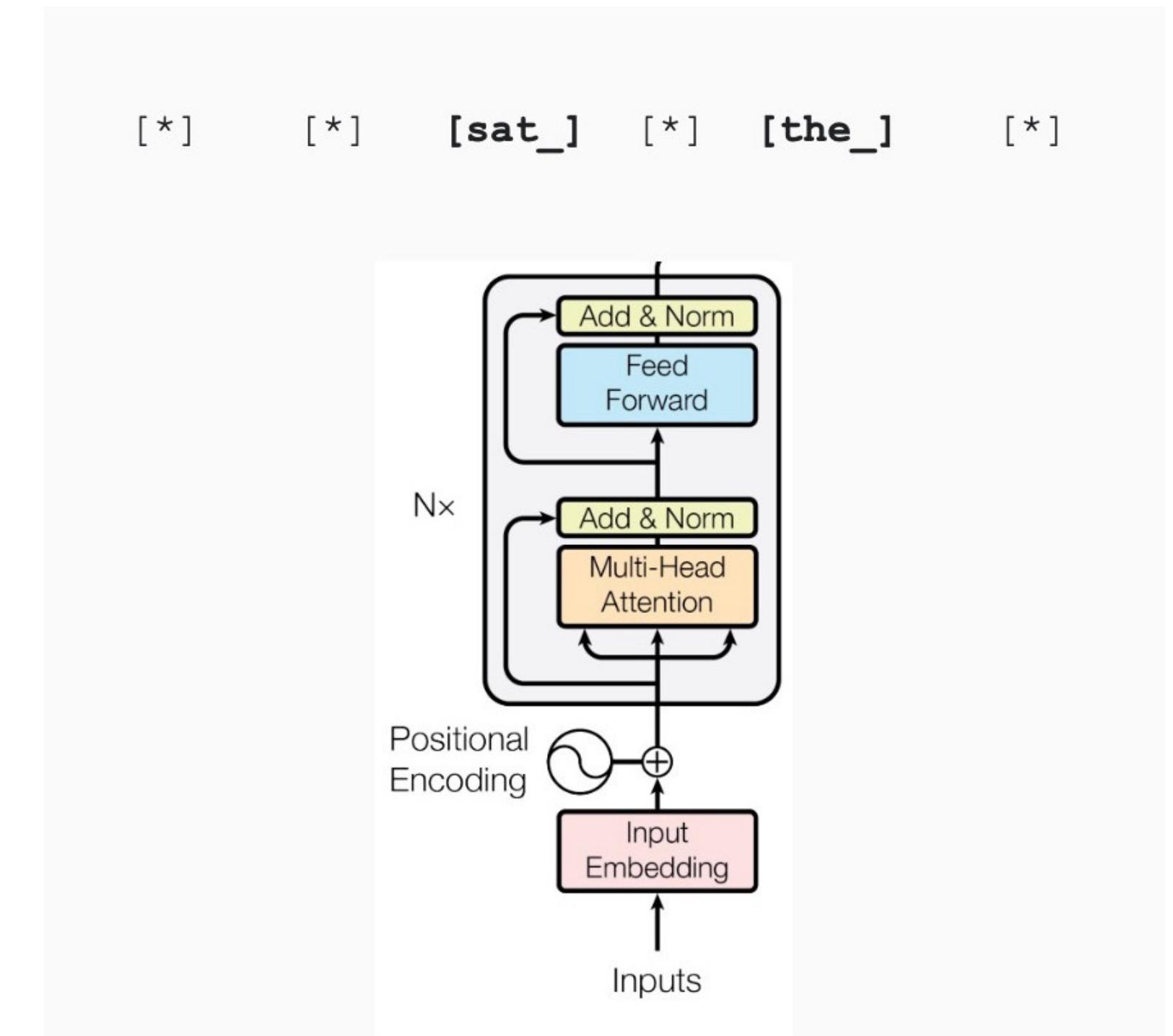
Translate EN-DE: This is good.
Summarize: state authorities dispatched...
Is this toxic: You look beautiful today!

[START] [The_] [cat_]

Transformer image source: "Attention Is All You Need" paper

BERT (2019)

- *Bidirectional* Encoder Representations from Transformers
- Encoder-only (no attention masking)
- 110M params
- 15% of all words masked out
- Was great, now dated



[The_] [cat_] [MASK] [on_] [MASK] [mat_]

From Lukas Beyer's lecture Intro to Transformers.

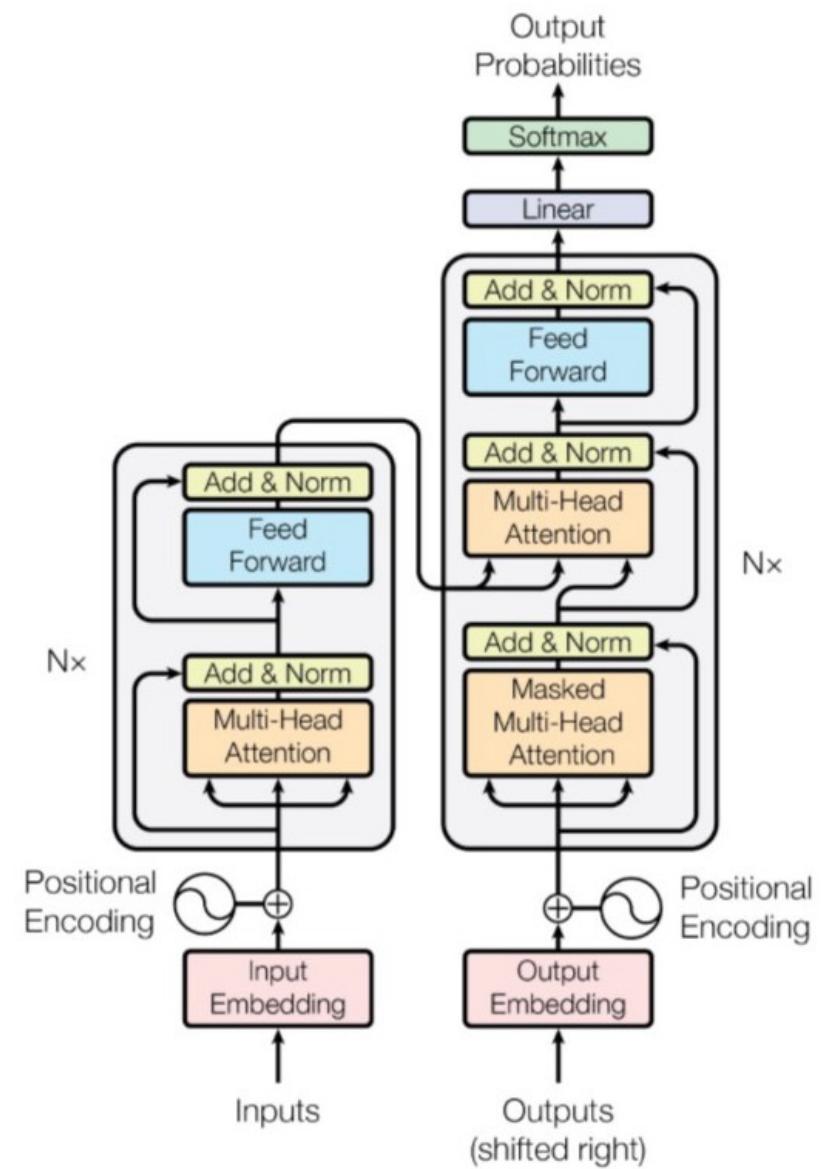
T5: Text-to-Text Transfer Transformer (2020)

- Input and output are both text strings
- Encoder-Decoder architecture
- 11B parameters
- Still could be a good choice for fine-tuning!

Das ist gut.

A storm in Attala caused 6 victims.

This is not toxic.



Translate EN-DE: This is good.

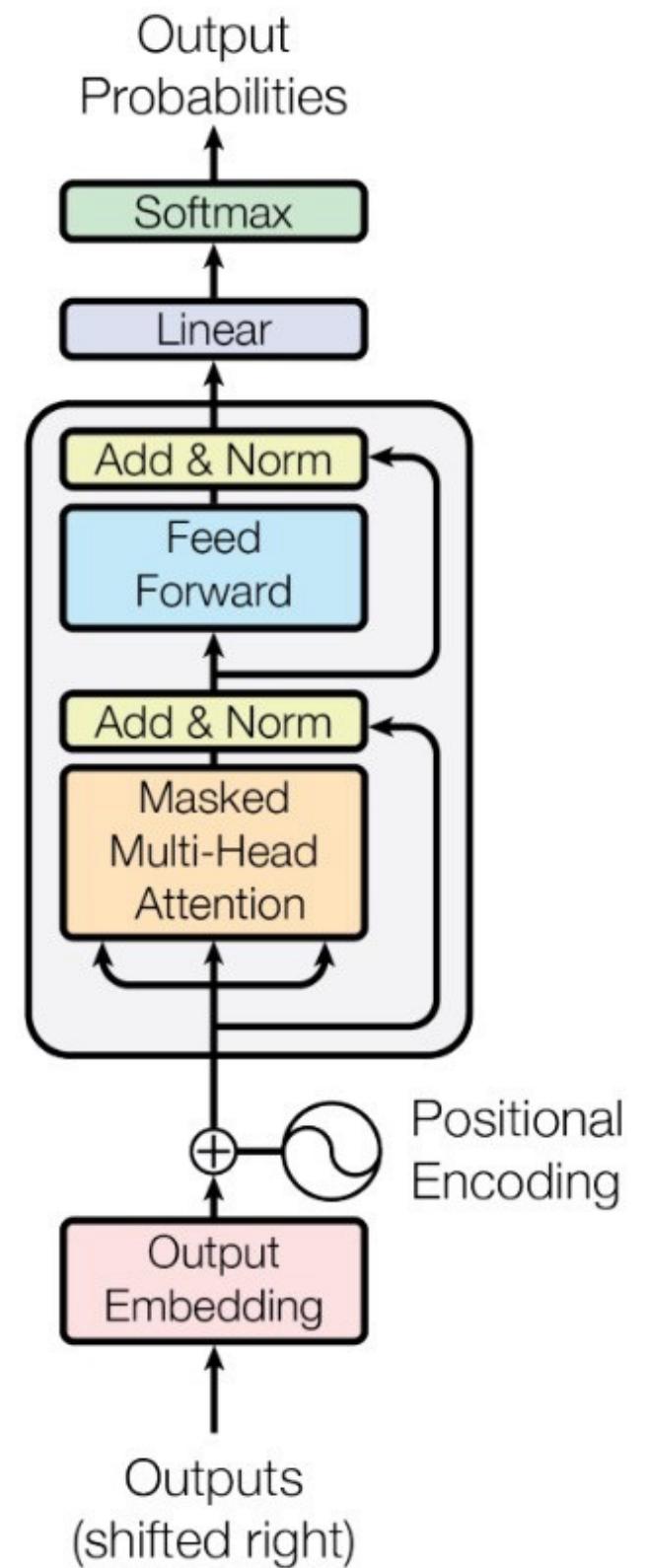
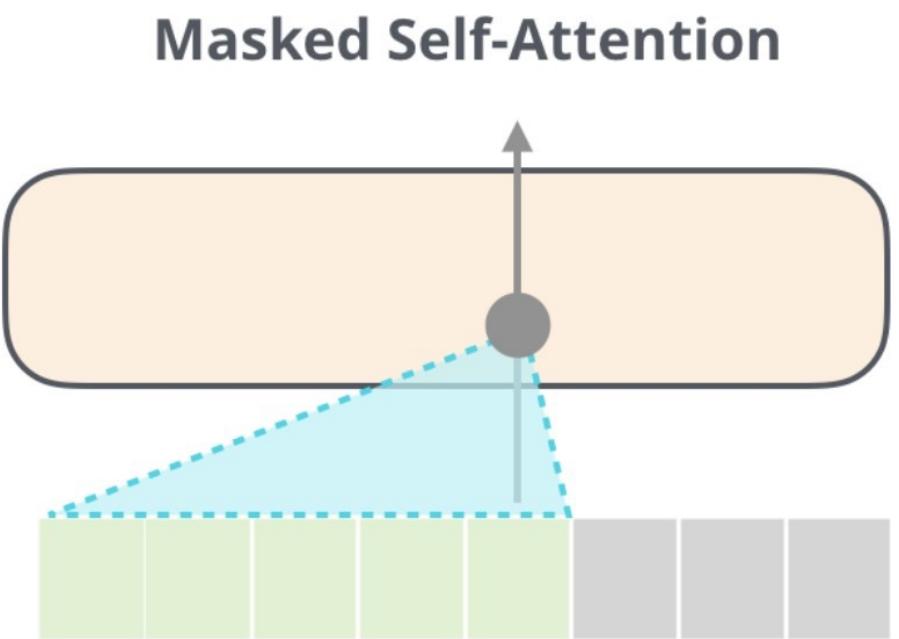
Summarize: state authorities dispatched...

Is this toxic: You look beautiful today!

GPT / GPT-2 (2019)

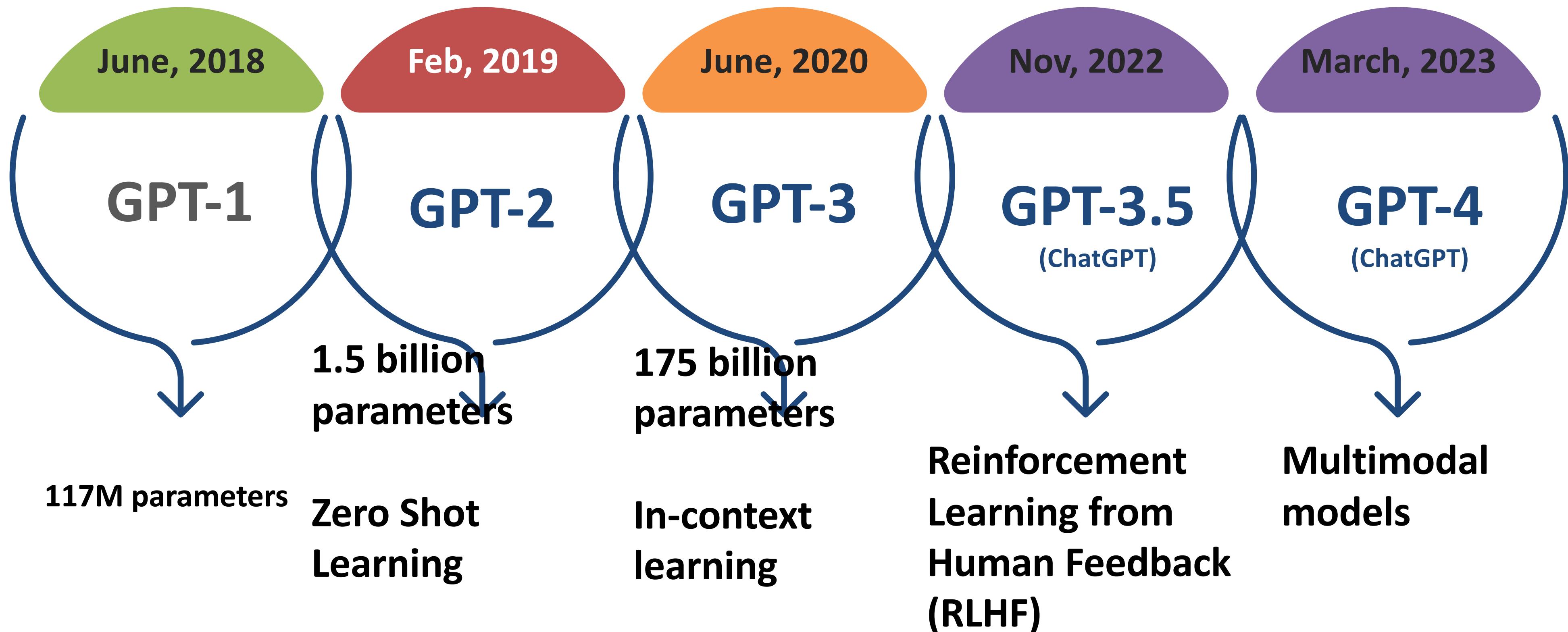
[sat_]

- Generative Pre-trained Transformer
- Decoder-only (uses masked self-attention)
- 1.5B



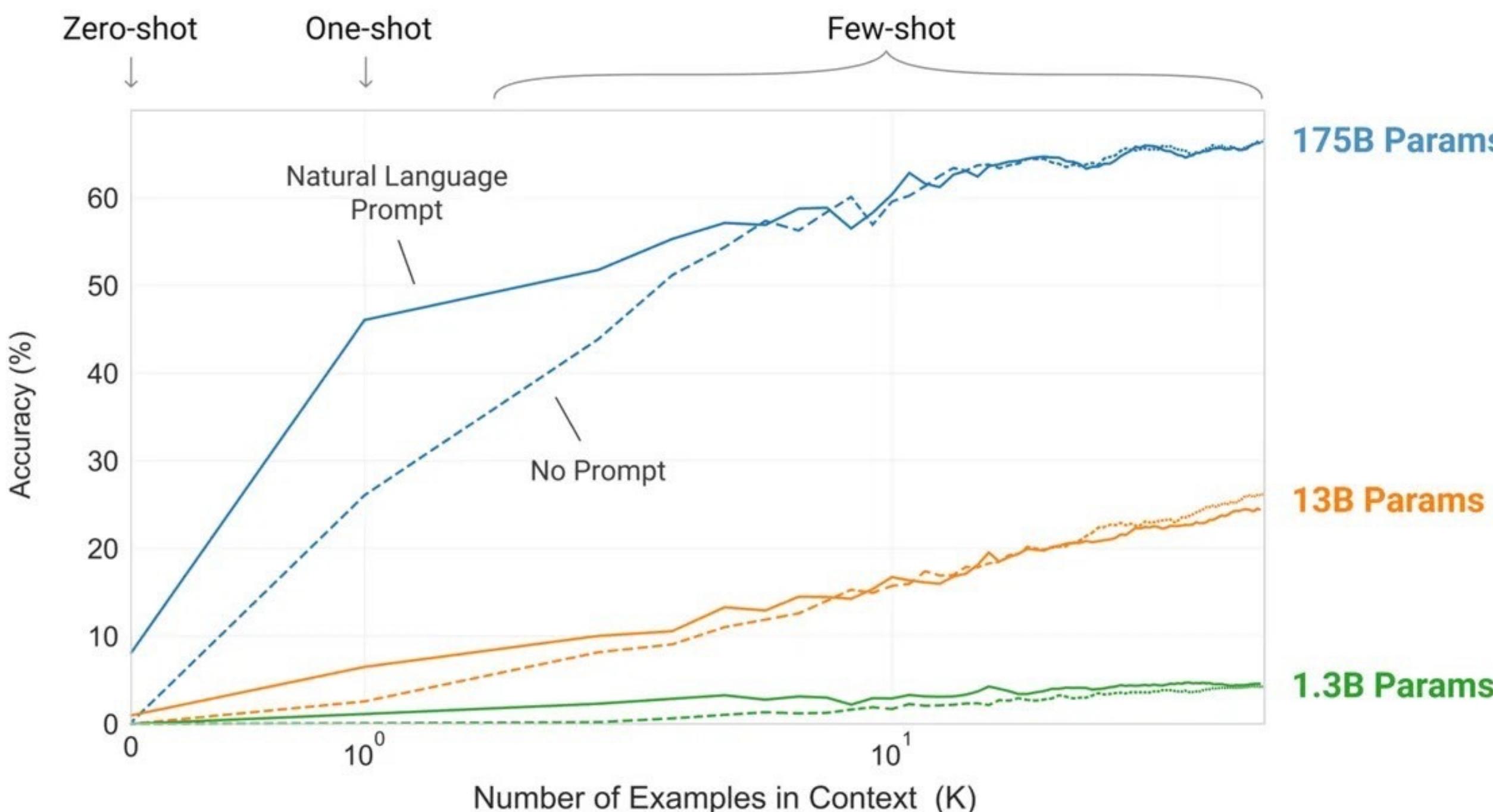
[START] [The_] [cat_]

GPT Timeline



GPT-3 (2020)

- Just like GPT-2, but 100x larger (175B params)
- Exhibited unprecedented few-shot and zero-shot learning



Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.

1 Translate English to French: ← task description
2 cheese => ← prompt

One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.

1 Translate English to French: ← task description
2 sea otter => loutre de mer ← example
3 cheese => ← prompt

Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

1 Translate English to French: ← task description
2 sea otter => loutre de mer ← examples
3 peppermint => menthe poivrée
4 plush girafe => girafe peluche
5 cheese => ← prompt

GPT-3 Training Data

- For a total of 500B tokens
- But trained on only 300B!

- Wikipedia (facts) (3.49%)
- Books1/BookCorpus (Smashwords) (7.8%)
- Books2 (Libgen or similar) (8.1%)
- WebText (Reddit links) (18.86%)
- Common Crawl (www) (61.75%)

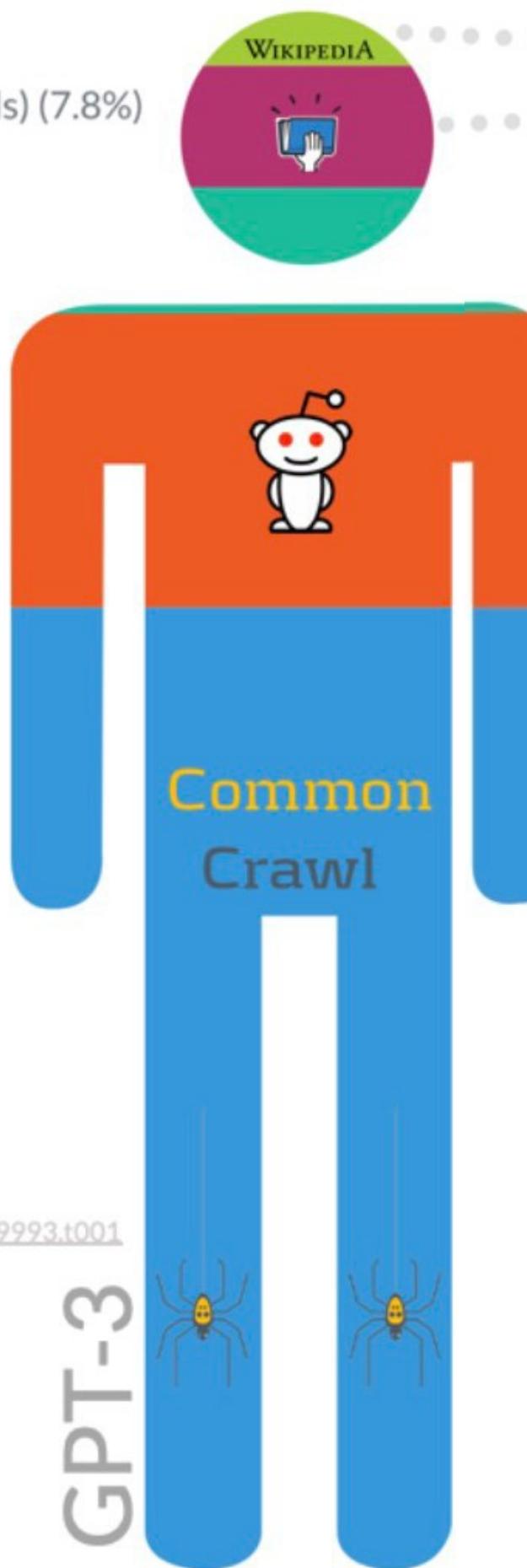
- Not to scale.
- Effective size by weighting (as % of total).
- Deduplication has been considered for Wikipedia.

Sources:
GPT3: <https://arxiv.org/abs/2005.14165>
The Pile v1: <https://arxiv.org/abs/2101.00027>
C4: <https://arxiv.org/abs/2104.08758>
Domains: <https://doi.org/10.1371/journal.pone.0249993.t001>

Alan D. Thompson. July 2021.
<https://lifearchitect.com.au/ai/>



LifeArchitect.ai/models



WebText (Reddit Submission Corpus)

HuffPost (news)
The New York Times (news)
BBC (news)
Twitter (discussion)
The Guardian (news)
The Washington Post (news)
and 4.3M+ more domains...

Common Crawl

(C4, cleaned/filtered, sorted by most tokens)

Google Patents (papers)
The New York Times (news)
Los Angeles Times (news)
The Guardian (news)
PLoS - Public Library of Science (papers)
Forbes (news)
HuffPost (news)
Patents.com - dead link (papers)
Scribd (books)
The Washington Post (news)
The Motley Fool (opinion)

InterPlanetary File System (mix)
Frontiers Media (papers)
Business Insider (news)
Chicago Tribune (news)
Booking.com (discussion)
The Atlantic (news)

Springer Link (papers)

Al Jazeera (news)

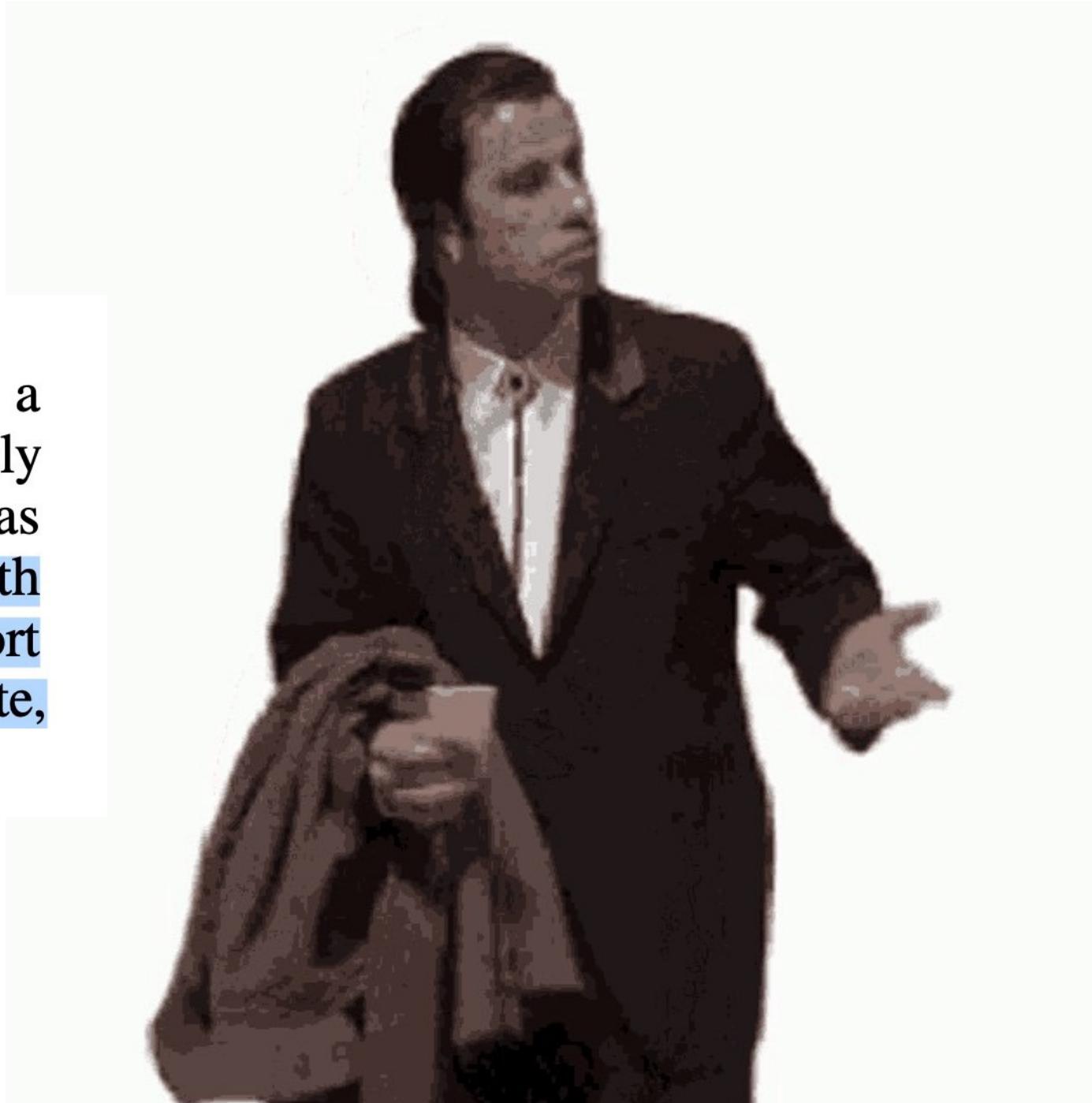
Kickstarter (discussion)

FindLaw Caselaw (papers)

National Center for Biotech Info (papers)
NPR (news)
and 90.9M+ more domains...

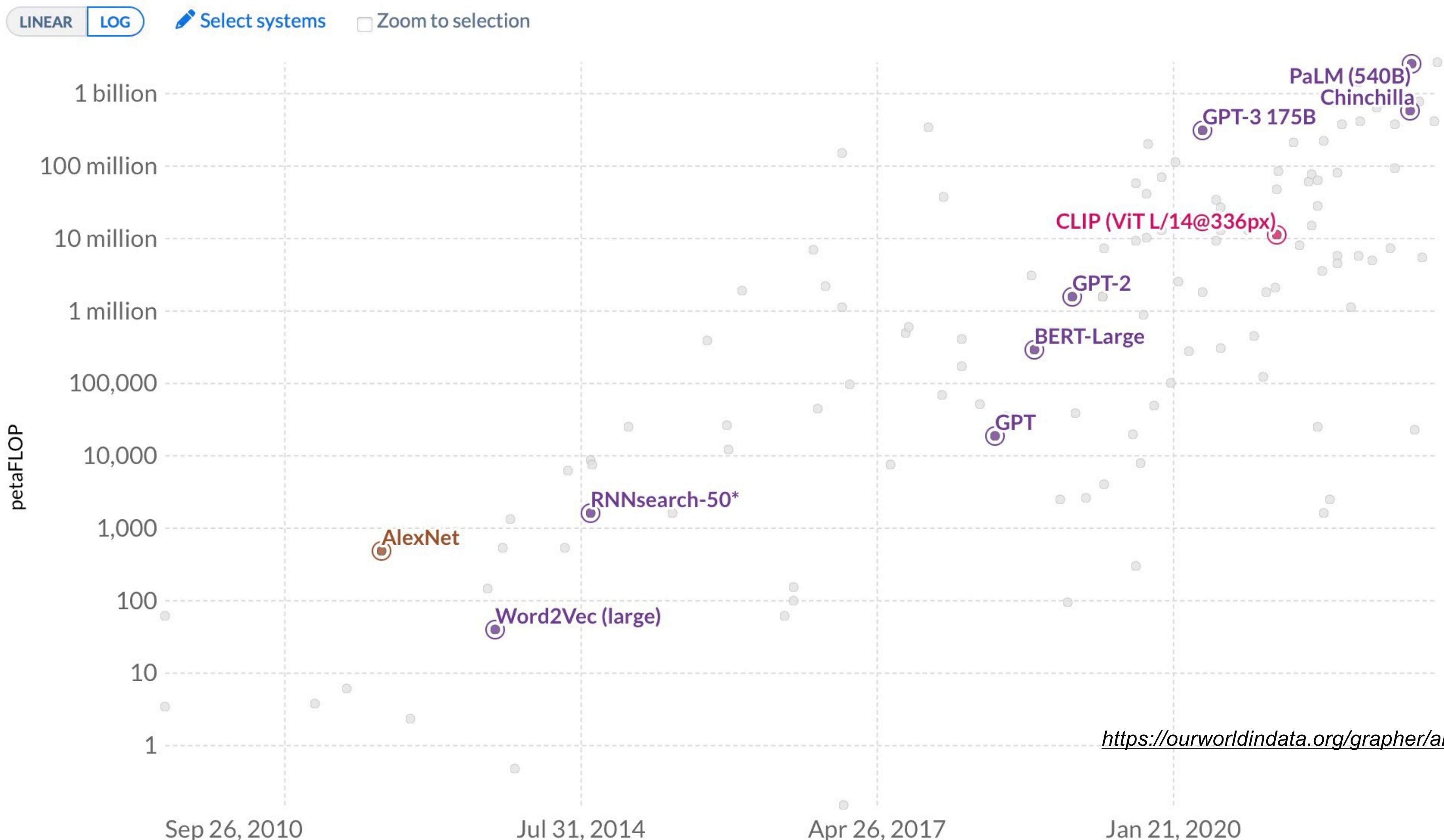
GPT-4 (2023)

This report focuses on the capabilities, limitations, and safety properties of GPT-4. GPT-4 is a Transformer-style model [39] pre-trained to predict the next token in a document, using both publicly available data (such as internet data) and data licensed from third-party providers. The model was then fine-tuned using Reinforcement Learning from Human Feedback (RLHF) [40]. Given both the competitive landscape and the safety implications of large-scale models like GPT-4, this report contains no further details about the architecture (including model size), hardware, training compute, dataset construction, training method, or similar.



Computation used to train notable AI systems

Computation is measured in petaFLOP, which is 10^{15} floating-point operations.

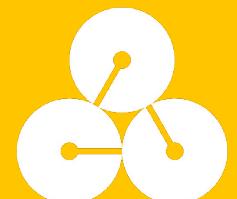


But what exactly is the relationship between
model size and dataset size?

But what exactly is the relationship between model size and dataset size?

Larger models tend to perform better, but the gains from increasing model size diminish unless the dataset size also increases.

Introducion to LLM APIs



Few LLM APIs

- OpenAI APIs
 - Models available: GPT-3.5 Turbo, GPT-4
 - Tasks: Completion, fine-tuning, function calling
- Anthropic
 - Models available: Claude Instant, Claude 2
 - Tasks: Completion.
- AWS Bedrock
 - Models available: Titan, Llama 2, and many more

OpenAI ChatGPT API calls

- Install OpenAI library
- Get an OpenAI API key

API calls are as easy as invoking a REST API

```
secret_key = 'copy and paste your key here'

import openai

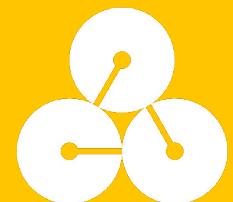
# Set up your OpenAI API credentials
openai.api_key = secret_key

# output = openai.Model.list()

output = openai.Completion.create(
    model="text-curie-001",
    prompt="tell me a joke",
)

print(output)
```

Introduction to Running LLMs Locally



Running LLMs Locally

- There are open-source ecosystems offering collections of chatbot LLMs similar to ChatGPT, which can be run locally on your computer.
- Why?
 - Offline Mode
 - Privacy and Security
 - Cost

GPT4All

- GPT4All is an LLM that you can install on your computer
- Try different models locally with UI setting

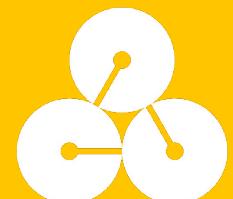
What models are supported by the GPT4All ecosystem?

Currently, there are six different model architectures that are supported:

1. GPT-J - Based off of the GPT-J architecture with examples found [here](#)
2. LLaMA - Based off of the LLaMA architecture with examples found [here](#)
3. MPT - Based off of Mosaic ML's MPT architecture with examples found [here](#)
4. Replit - Based off of Replit Inc.'s Replit architecture with examples found [here](#)
5. Falcon - Based off of TII's Falcon architecture with examples found [here](#)
6. StarCoder - Based off of BigCode's StarCoder architecture with examples found [here](#)

**Start With
GPT-J**

Libraries to Develop LLM Applications



Popular libraries for LLM Ecosystem

- [LangChain](#)
 - Popular Python library for chaining multiple generative models together
- [Streamlit](#)
 - For building ChatGPT like web interface
- [Hugging Face](#)
 - Python tools for pre-processing, training, fine-tuning, and deployment
- [Pinecone](#), [Chroma DB](#), [Milvus](#)
 - Vector databases

AI Feature Proposal

Due 25th February

You will need to integrate an AI-based feature into your web application. Before you begin your implementation, you will need to submit a proposal for the feature you will be implementing, thinking about the scope and the overall design of the feature, and how you will integrate AI. This proposal does not need to be a precise design document but rather should be a rough sketch to make your idea a little more concrete.

Upon submitting your proposal, we will provide feedback on the above, and either return "approve" or "needs revision".