

LLM Bootcamp
2023

LLM Foundations
Sergey Karayev

APRIL 21, 2023



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 Running a Transformer

03

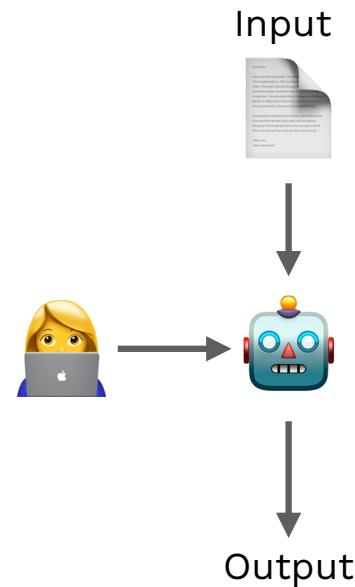
TRAINING & INFERENCE

00

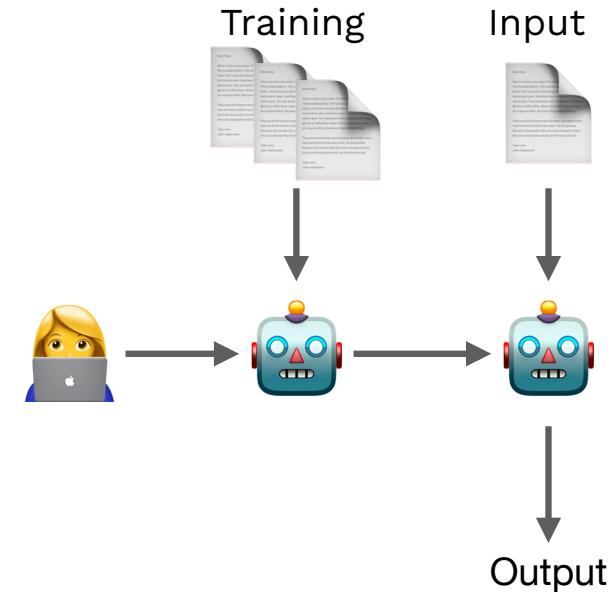
Foundations of Machine Learning



Traditional Programming vs Machine Learning



Software 1.0



Software 2.0

Types of Machine Learning

Unsupervised Learning

Learn structure of data to generate more data

"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



Converging on just...

Supervised or Self-supervised Learning

"This product does what it is supposed __" → "to."



→ *cat*



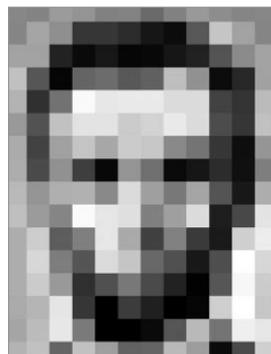
→ "Hey
Siri"



→ *next
move*

Inputs and outputs are always just numbers

What we see



Input

Output

"Lincoln"

What the machine "sees"

157	153	174	168	150	152	129	151	172	161	155
155	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	269	185	215	211	198	199	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	155	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

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

Why is this hard?

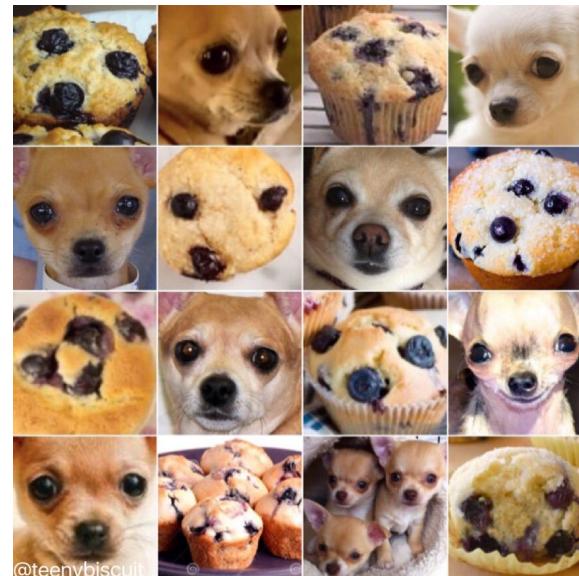
- Infinite variety of inputs can all mean the same thing
- Meaningful differences can be tiny
- Structure of the world is complex



"I loved this movie"

"As good as The Godfather"

"🔥 no cap"

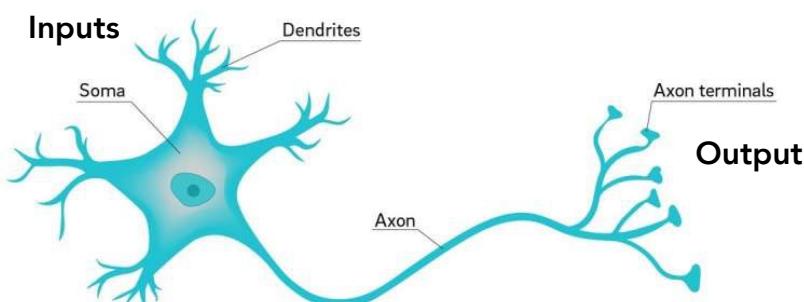
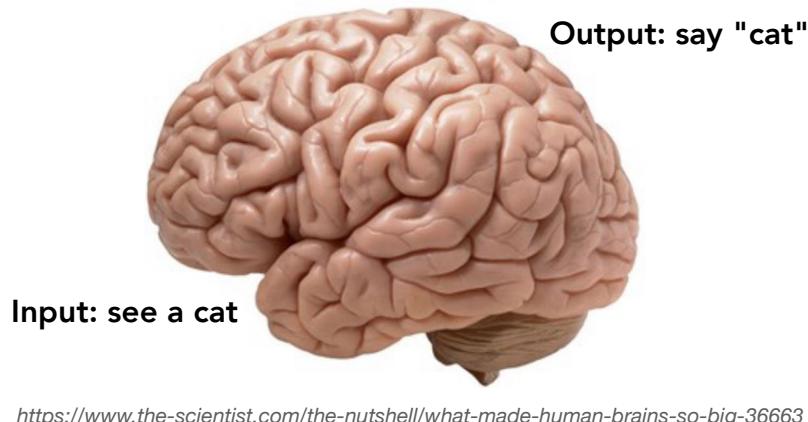


@teenybiscuit

How is it done?

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

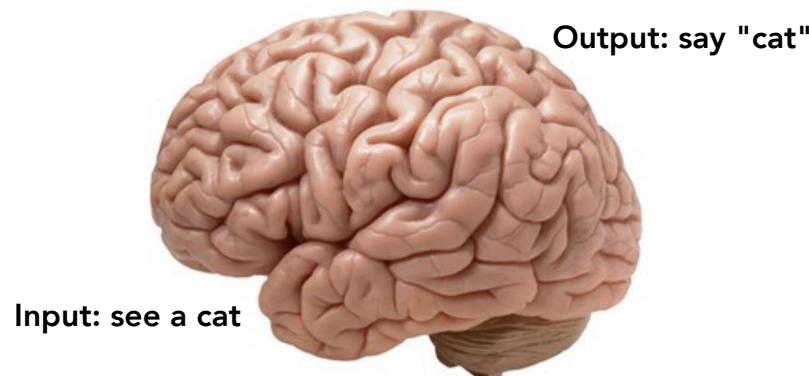
Inspiration



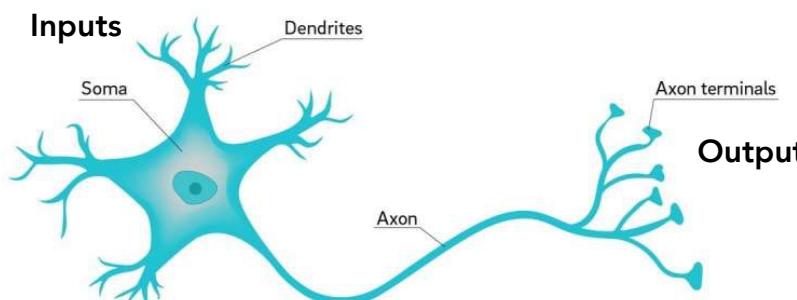
<https://medicalxpress.com/news/2018-07-neuron-axons-spindly-theyre-optimizing.html>

- 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

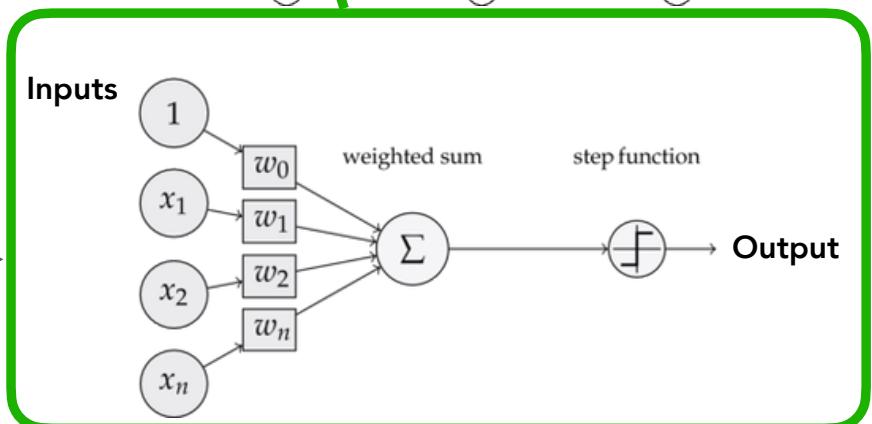
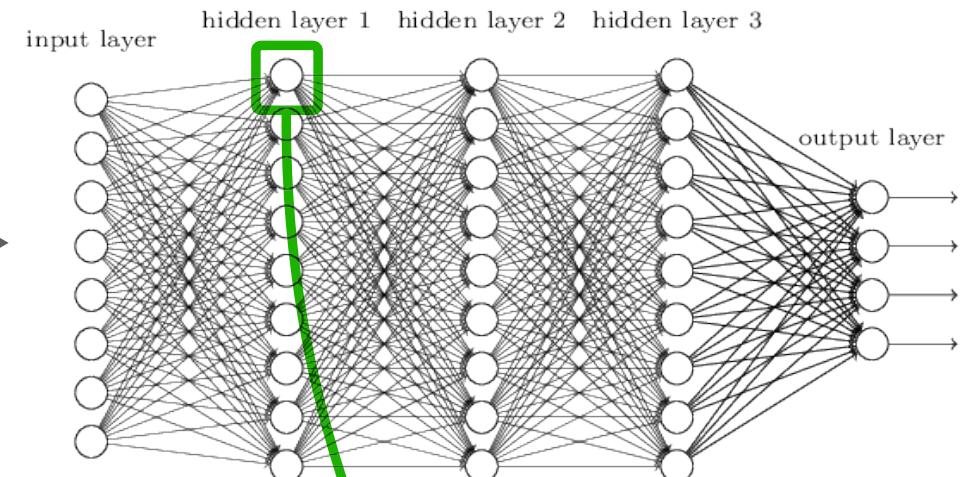
Formalization



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



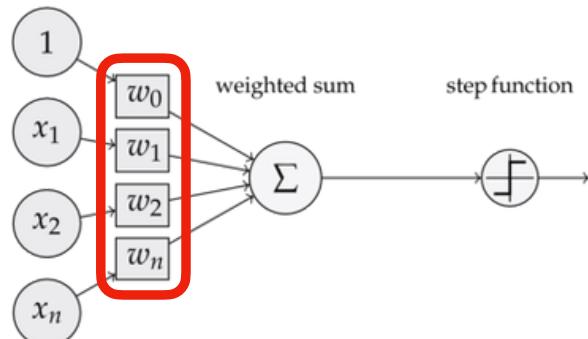
<https://medicalxpress.com/news/2018-07-neuron-axons-spindly-theyre-optimizing.html>



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

A "perceptron" is a vector of numbers

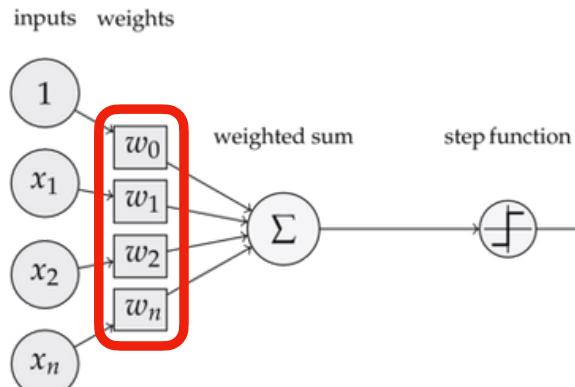
inputs weights



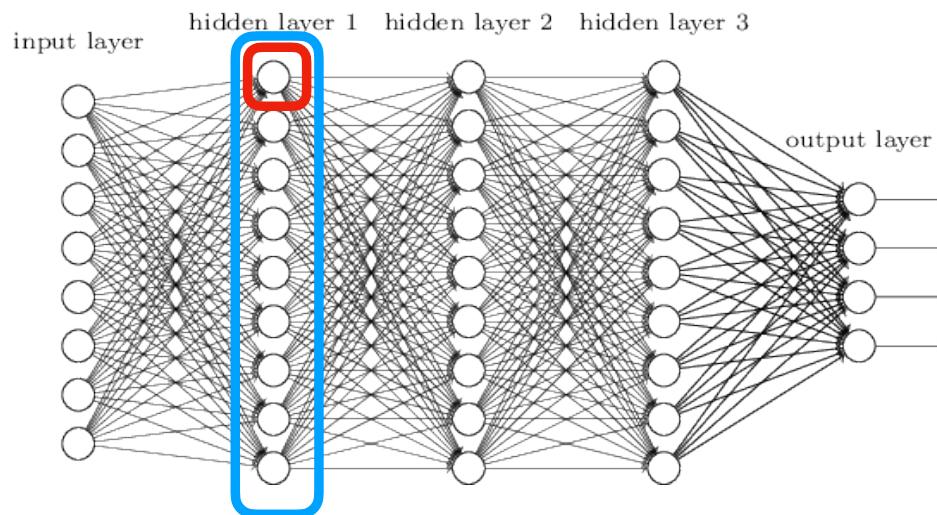
$$[0.58 \quad 0.841 \quad 0.835 \quad 0.18]$$

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

A "layer" is a matrix of numbers



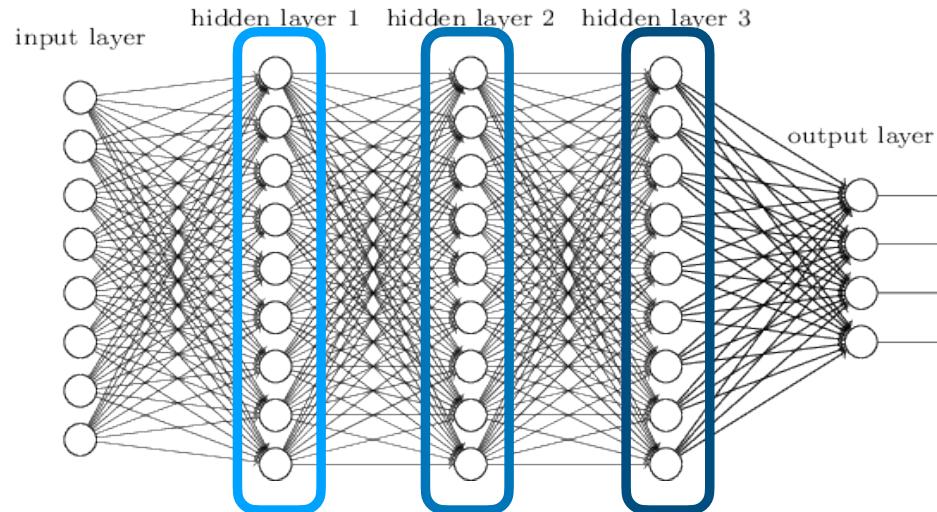
<https://www.jessicayung.com/explaining-tensorflow-code-for-a-multilayer-perceptron/>



[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"



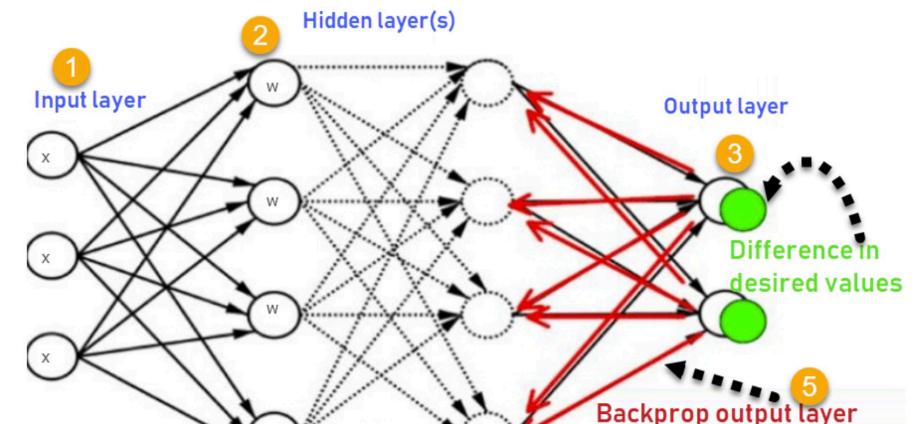
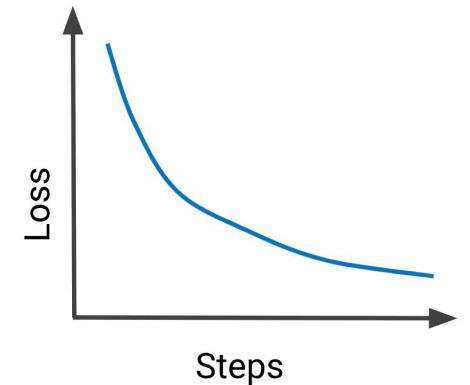
[[0.32 0.866 0.151 0.555]
[
[[0.58 0.841 0.835 0.18 1
[
[[0.32 0.866 0.151 0.555]
[
[[0.467 0.414 0.562 0.555]
[
[[0.186 0.937 0.131 0.448]
[
[[0.892 0.567 0.111 0.678]
[
[[0.483 0.478 0.424 0.844]
[
[[0.774 0.963 0.213 0.569]
[
[[0.014 0.509 0.21 0.26]
[
[[0.142 0.991 0.105 0.211]
[
[[0.789 0.242 0.774 0.019]]

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

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

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

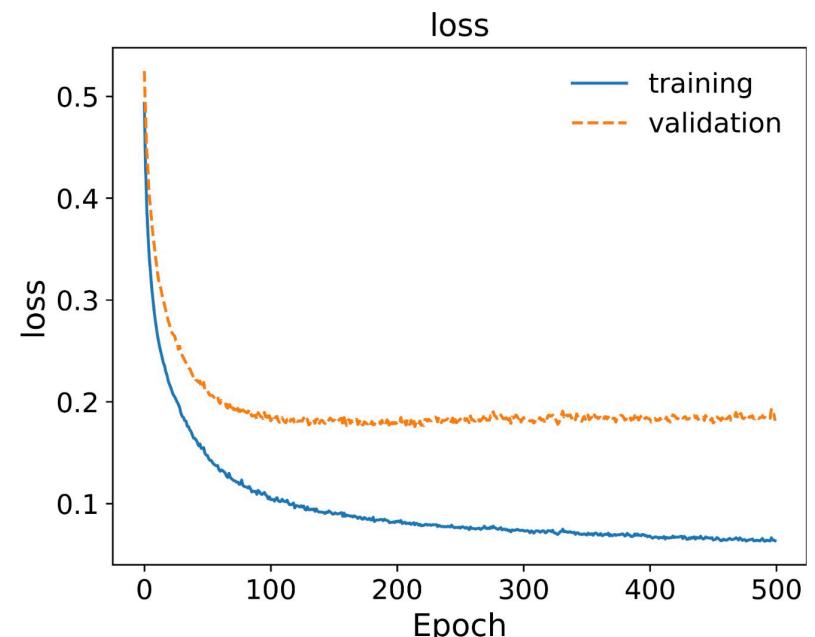


<https://www.guru99.com/backpropagation-neural-network.htm>

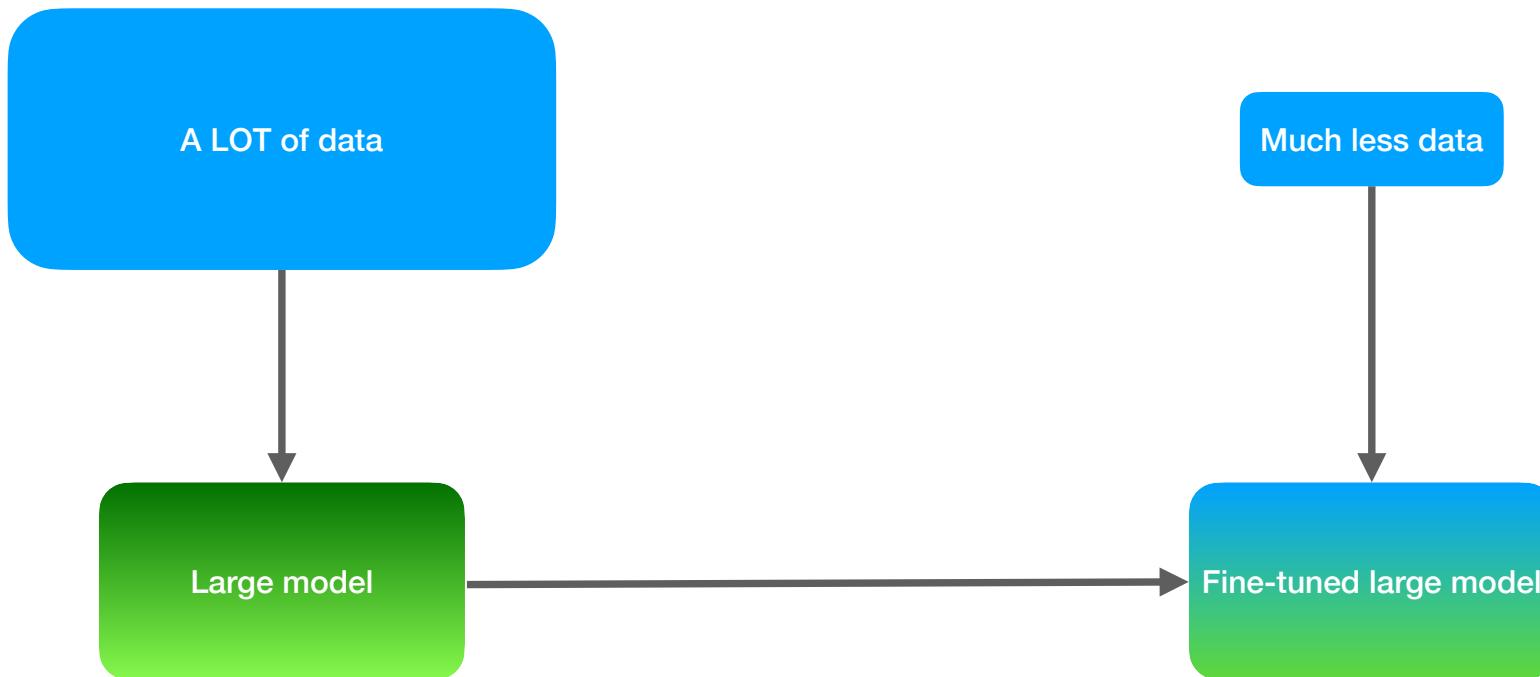
<https://developers.google.com/machine-learning/testing-debugging/metrics/interpretic>

Dataset Splitting

- Split (X, y) into training (~80%), validation (~10%), and test (~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!**

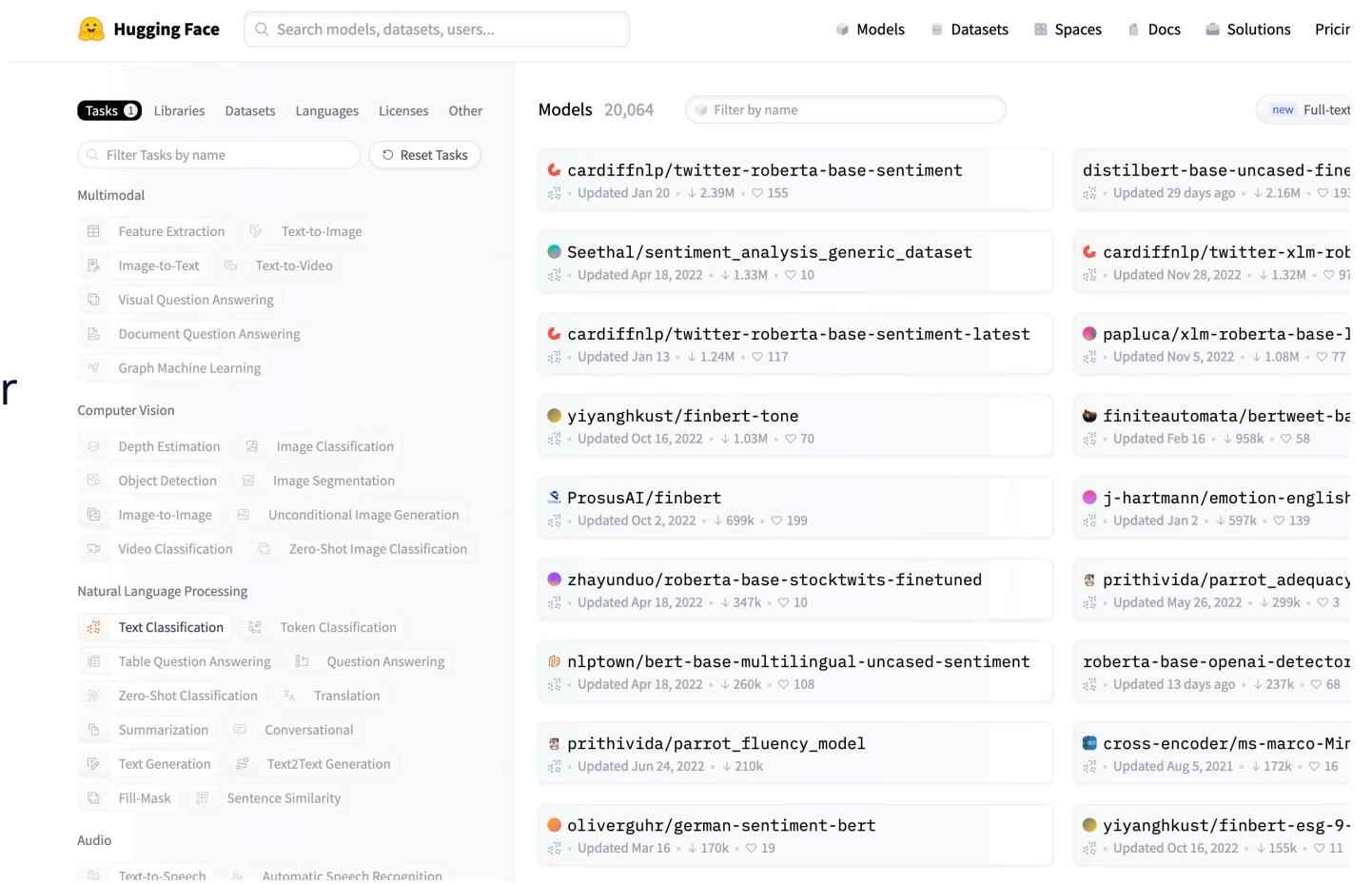


Pre-training:
slow training on a lot of data

Fine-tuning:
fast training on a little data

Model Hubs

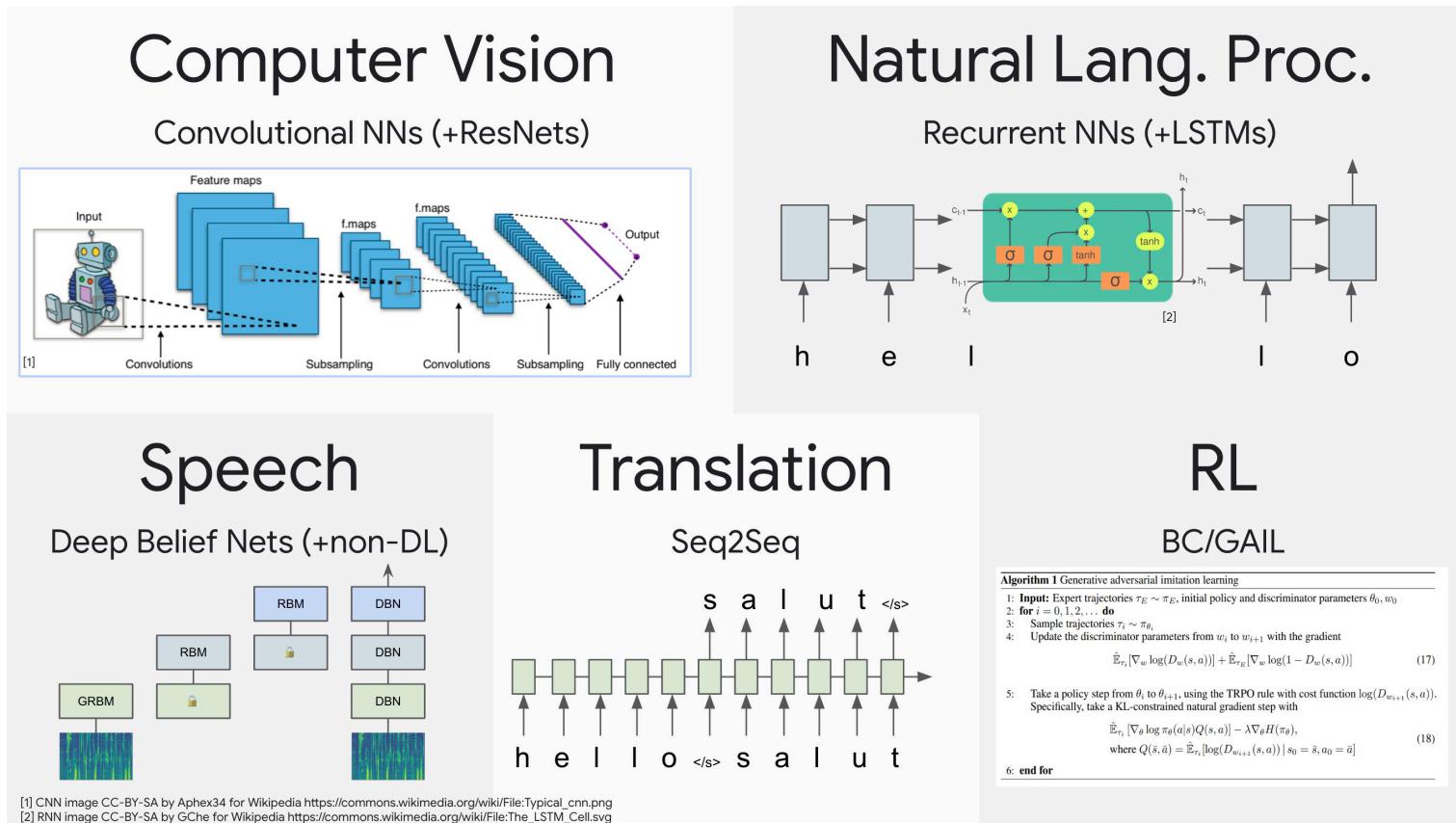
- People share pre-trained models!
- 😊: the most popular Model Hub
 - 180K models
 - 30K datasets



The screenshot shows the Hugging Face Model Hub interface. At the top, there's a search bar with the placeholder "Search models, datasets, users...". Below the search bar are navigation tabs: Models, Datasets, Spaces, Docs, Solutions, and Pricing. The "Models" tab is selected. On the left, there's a sidebar with a "Hugging Face" logo and a "Tasks" section containing "Libraries", "Datasets", "Languages", "Licenses", and "Other". Below "Tasks" is a search bar for "Filter Tasks by name" and a "Reset Tasks" button. The main content area is titled "Models 20,064" and features a "Filter by name" input field and a "new Full-text" link. The models are listed in a grid, each with a profile picture, the model name, the last update date, the number of downloads (indicated by a downward arrow), and the number of stars (indicated by a heart icon). The categories shown in the sidebar are Multimodal, Computer Vision, Natural Language Processing, and Audio.

Model Name	Last Updated	Downloads	Stars
cardiffnlp/twitter-roberta-base-sentiment	Updated Jan 20	2.39M	155
distilbert-base-uncased-fine	Updated 29 days ago	2.16M	191
Seethal/sentiment_analysis_generic_dataset	Updated Apr 18, 2022	1.33M	10
cardiffnlp/twitter-roberta-base-sentiment-latest	Updated Jan 13	1.24M	117
yiyangkust/finbert-tone	Updated Oct 16, 2022	1.03M	70
ProsusAI/finbert	Updated Oct 2, 2022	699k	199
zhayunduo/roberta-base-stocktwits-finetuned	Updated Apr 18, 2022	347k	10
nlptown/bert-base-multilingual-uncased-sentiment	Updated Apr 18, 2022	260k	108
prithivida/parrot_adequacy	Updated May 26, 2022	299k	3
roberta-base-openai-detector	Updated 13 days ago	237k	68
cross-encoder/ms-marco-Mr	Updated Aug 5, 2021	172k	16
oliverguhr/german-sentiment-bert	Updated Mar 16	170k	19
yiyangkust/finbert-esg-9-	Updated Oct 16, 2022	155k	11

Before ~2020: each task had its own NN architecture



Now: all is Transformers



Transformer cartoon (DALL-E)

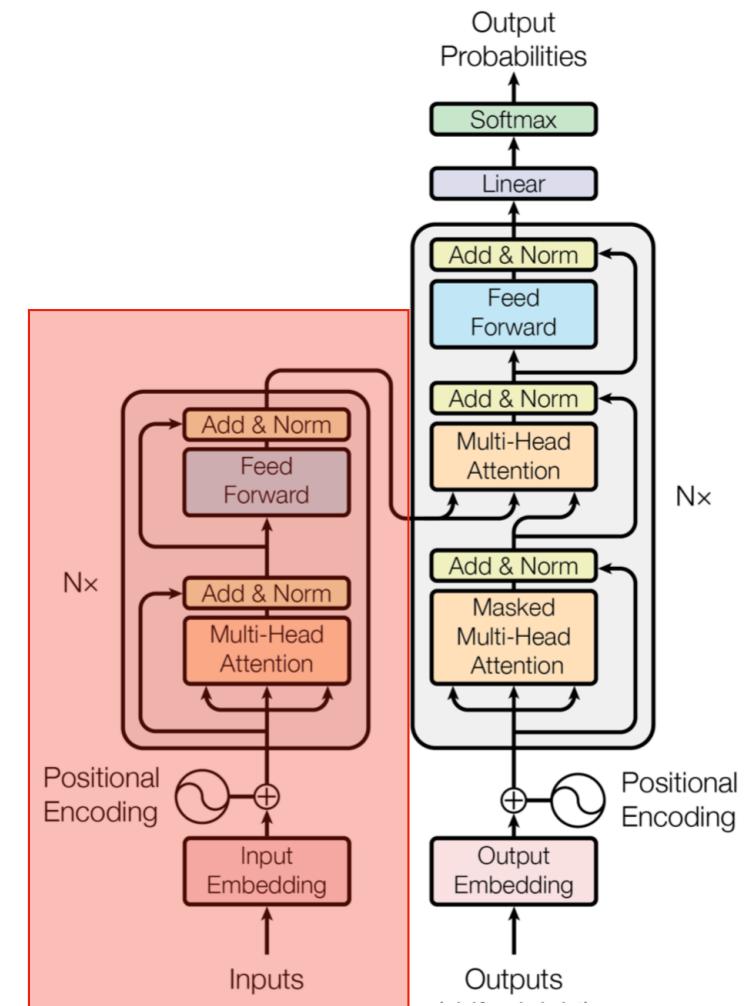
01

The Transformer Architecture



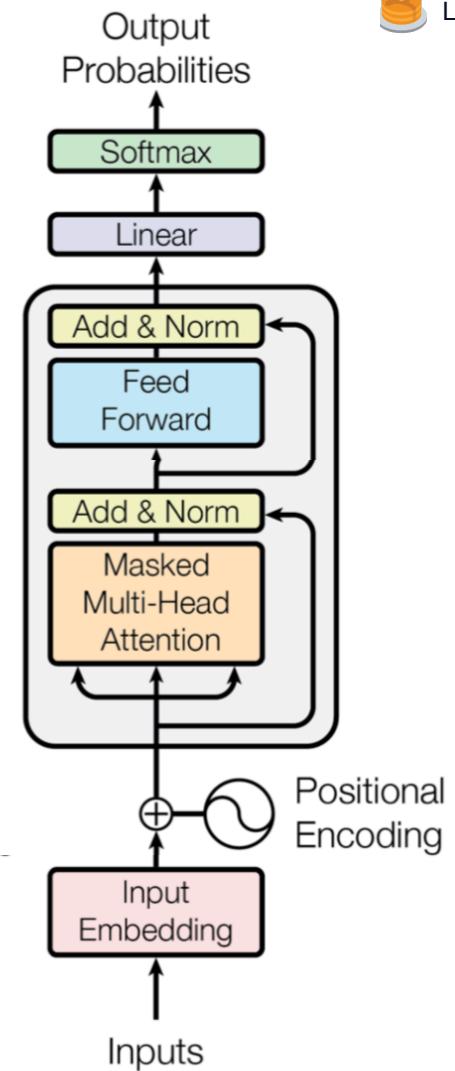
Attention is all you need (2017)

- 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



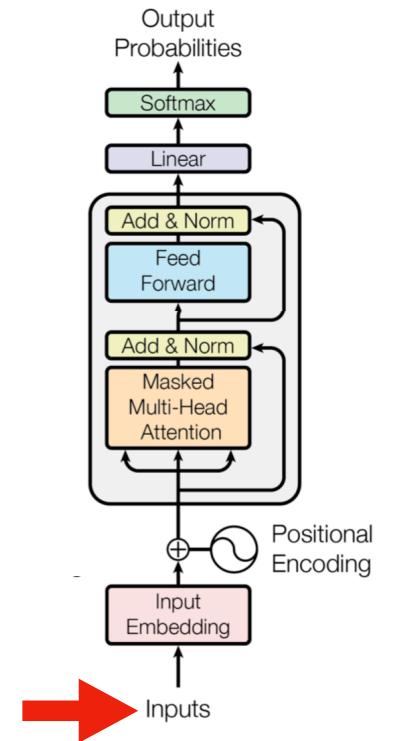
Transformer Decoder Overview

- 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.



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, sundress, ., <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. 3 -> [0, 0, 0, 1, 0, 0, 0, 0, ...]



Input Embedding

- 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"
- Solution: learn an embedding matrix!
 - (The simplest NN layer type)

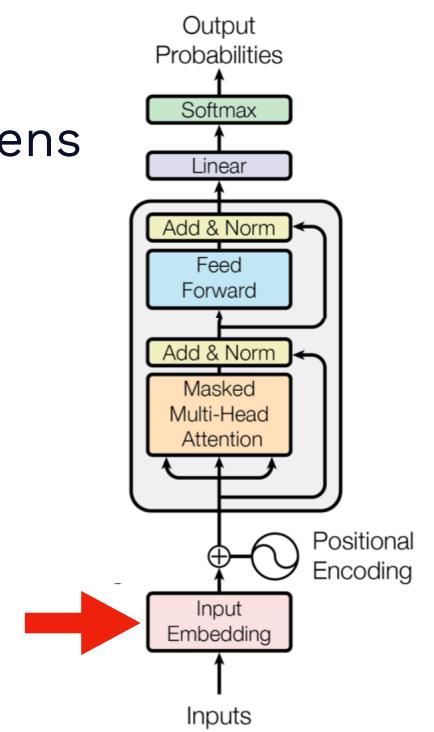
	1	0	0	...	0	0	0
aardvark	1	0	0	...	0	0	0
black	0	0	...	1	...	0	0
cat	0	0	...	1	...	0	0
duvet	0	0	...	1	...	0	0
zombie	0	0	0	...	0	0	1

$V \times V$

$V \times E$
*embedding
matrix*

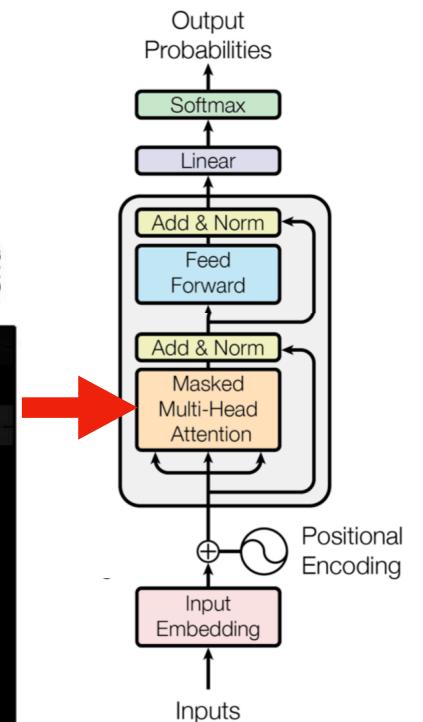
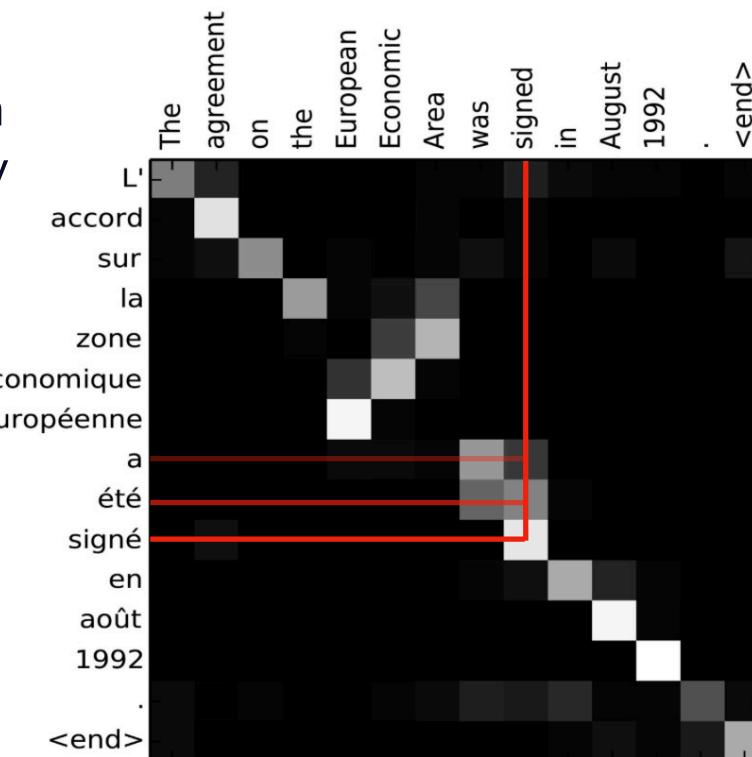
aardvark	0.97	0.03	0.15	0.04
black	0.07	0.01	0.20	0.95
cat	0.98	0.98	0.45	0.35
duvet	0.01	0.84	0.12	0.02
zombie	0.74	0.05	0.98	0.93

$V \times E$



Attention

- (Ignore "Masked Multi-Head" for now)
- Key insight: for a given token in the output sequence, only one or a few tokens in the input sequence are most important
- Introduced in 2015 for translation tasks



<https://lilianweng.github.io/posts/2018-06-24-attention/>

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

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

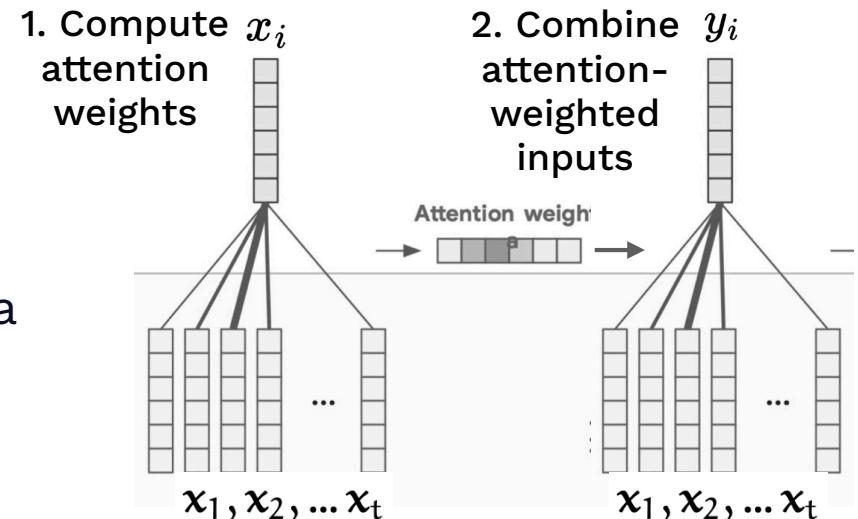
$$\mathbf{y}_i = \sum_j w_{ij} \mathbf{x}_j$$

- weight is just dot product between input vectors

$$w'_{ij} = \mathbf{x}_i^T \mathbf{x}_j$$

- (made to sum to 1)

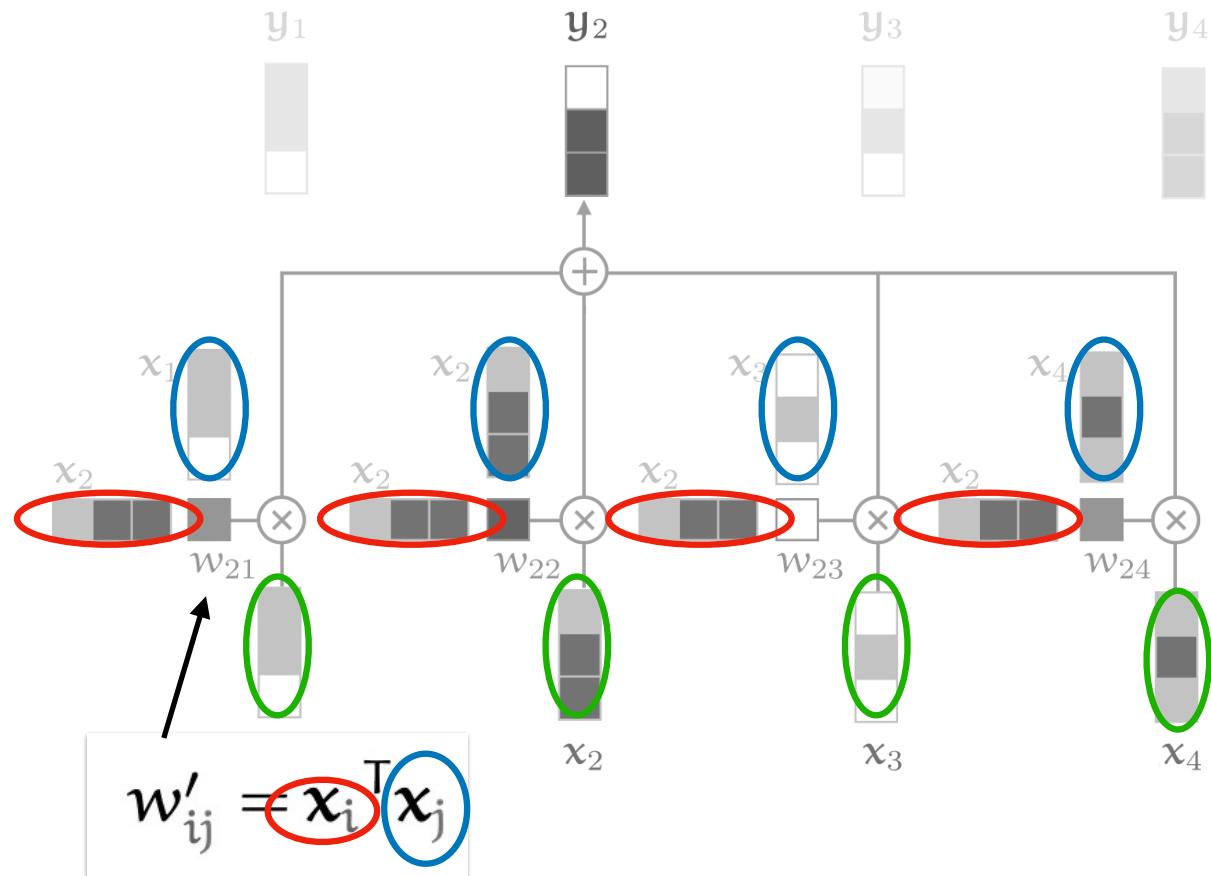
$$w_{ij} = \frac{\exp w'_{ij}}{\sum_j \exp w'_{ij}}$$



Basic self-attention

$$y_i = \sum_j w_{ij} x_j$$

- Note that every input vector x_i is used in 3 ways:
 - Query x_i
 - Key $x_1, x_2, \dots x_t$
 - Value $x_1, x_2, \dots x_t$



$$w'_{ij} = x_i^T x_j$$

Basic self-attention

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

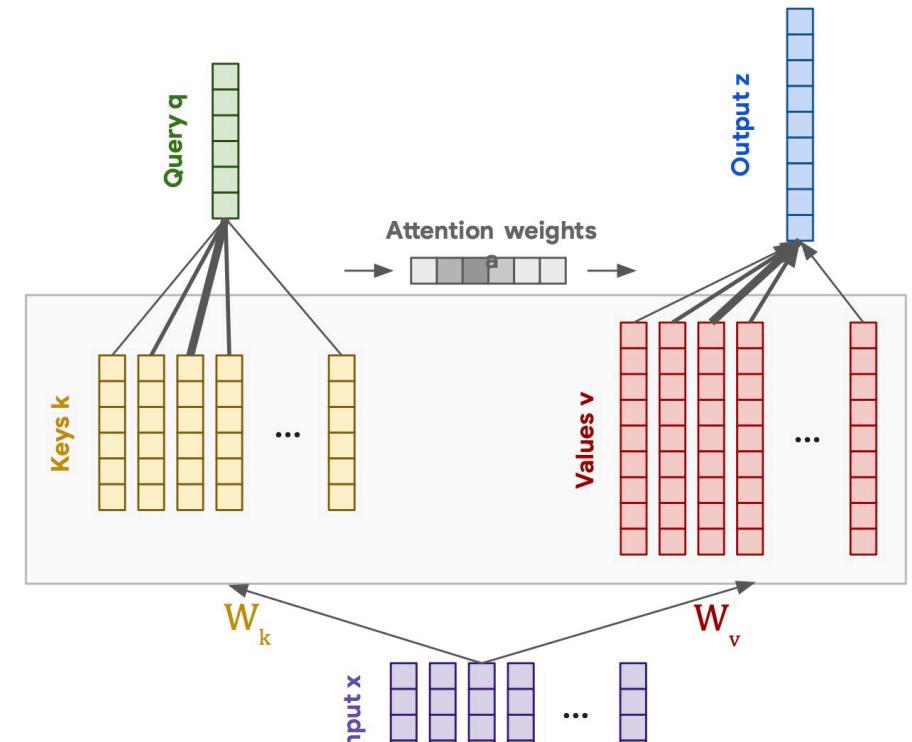
$$q_i = W_q x_i \quad k_i = W_k x_i \quad v_i = W_v x_i$$

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

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

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

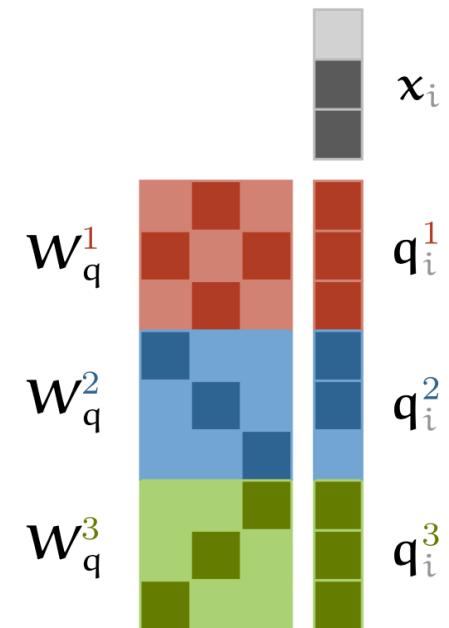
<http://www.peterbloem.nl/bloa/transformers>



<http://lucasbeyer.be/transformer>

Multi-head attention

- 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.
 - (Actually implemented as a single matrix, anyway.)

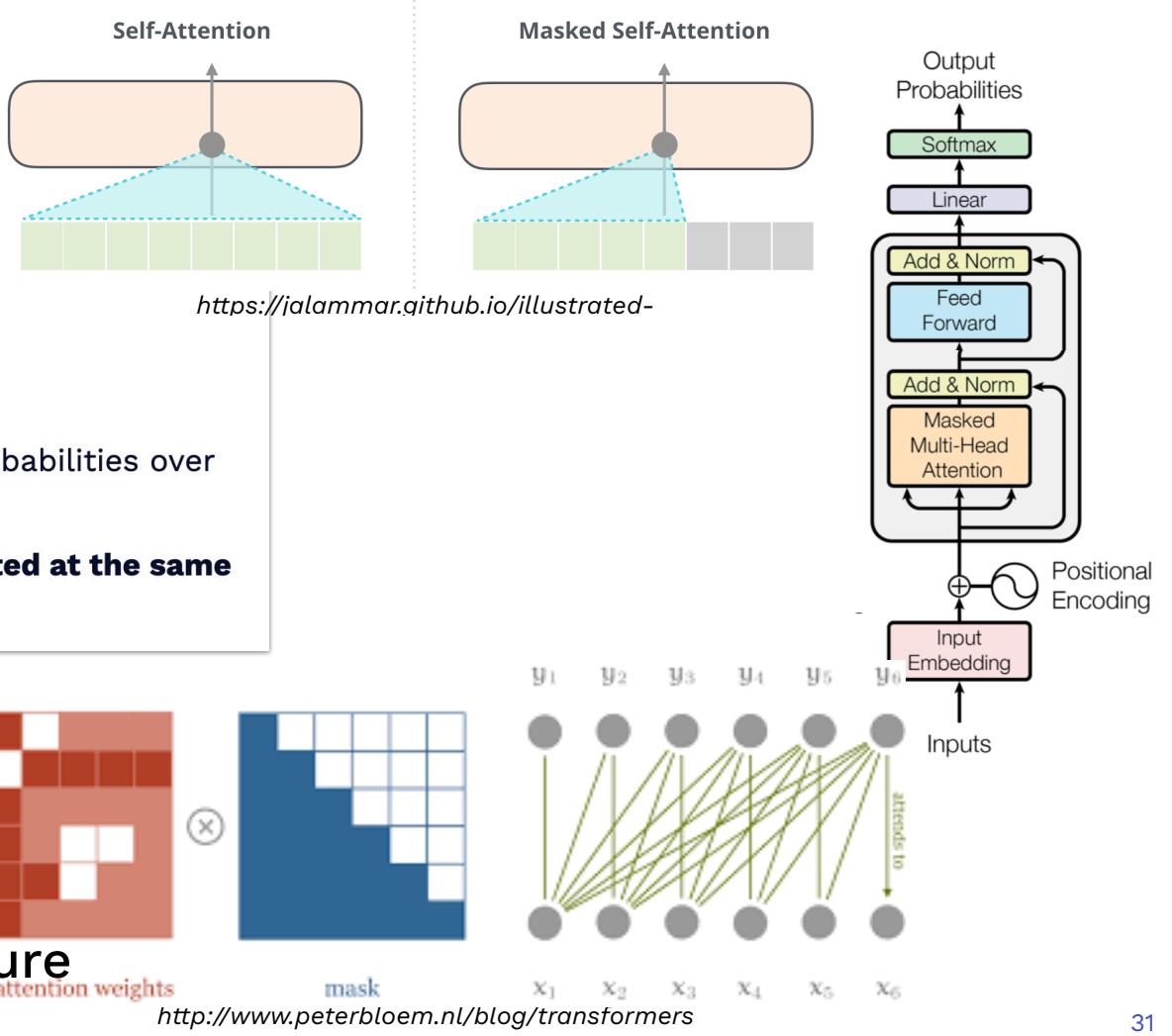


3-headed attention

Masking attention

In training:

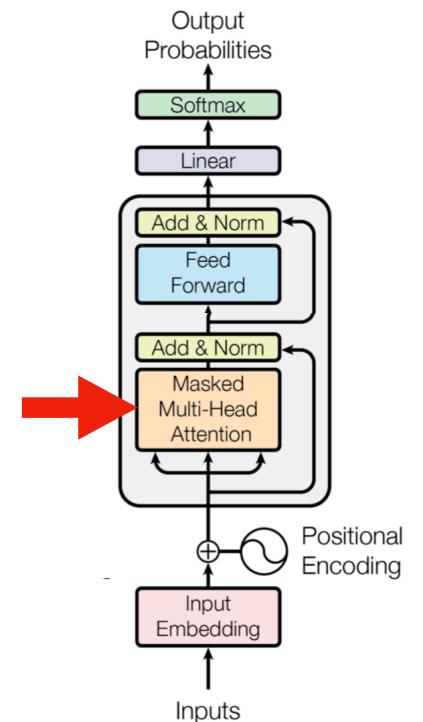
- Inputs: a sequence of N tokens
 - [It's, a, blue, <BLANK>, <BLANK>, .]
- Ground-truth Outputs: a sequence of N tokens
 - [a, blue, sundress, <BLANK>, <BLANK>, ...]
- Actual Outputs: Instead of words, vectors of probabilities over the vocabulary
 - **Crucially: all output probabilities are computed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



$$\mathbf{q}_i = \mathbf{W}_q \mathbf{x}_i$$

$$\mathbf{k}_i = \mathbf{W}_k \mathbf{x}_i$$

$$\mathbf{v}_i = \mathbf{W}_v \mathbf{x}_i$$



LLMBC 2023

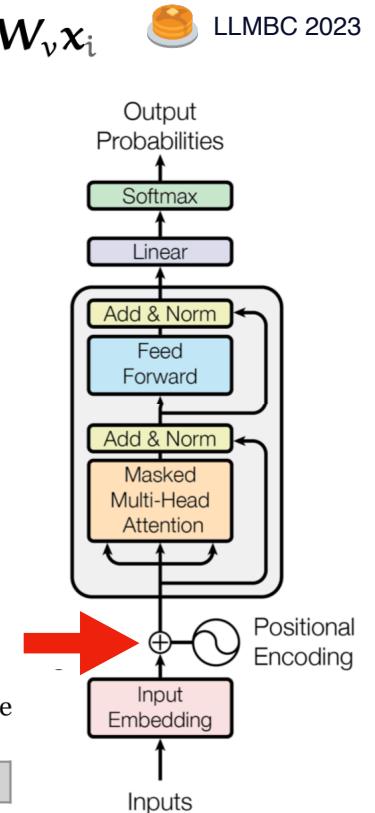
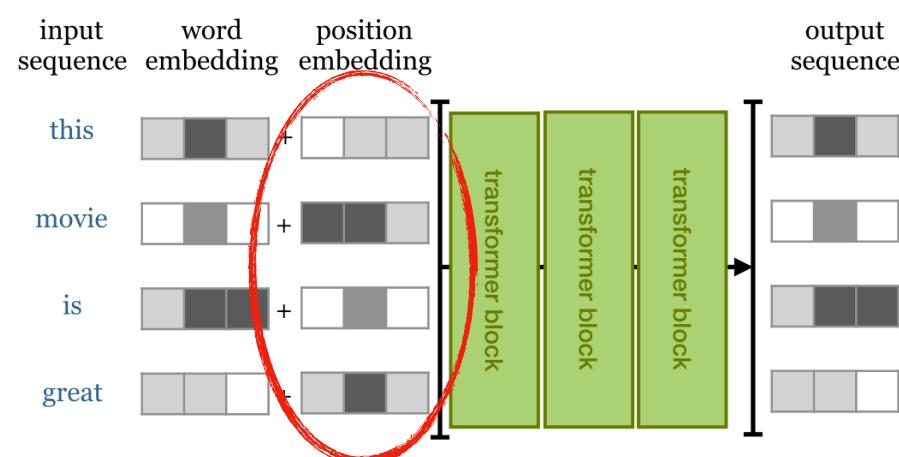
Positional Encoding

- 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

$$w'_{ij} = \mathbf{q}_i^\top \mathbf{k}_j$$

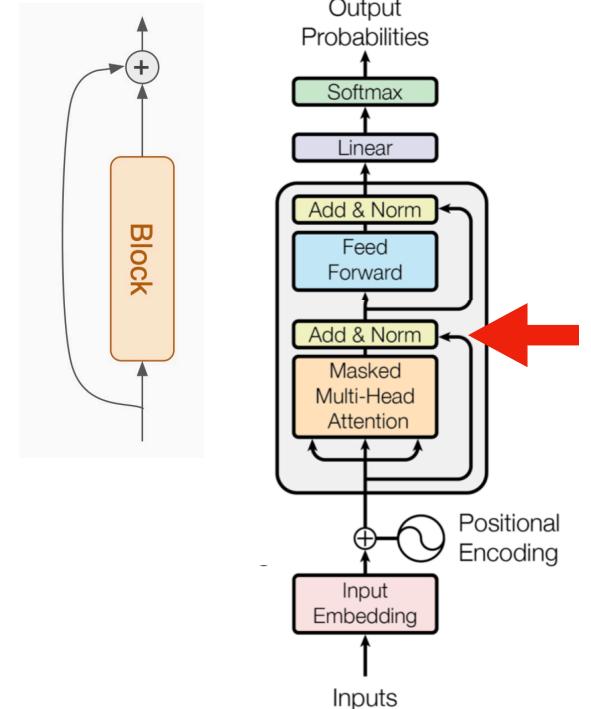
$$w_{ij} = \text{softmax}(w'_{ij})$$

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



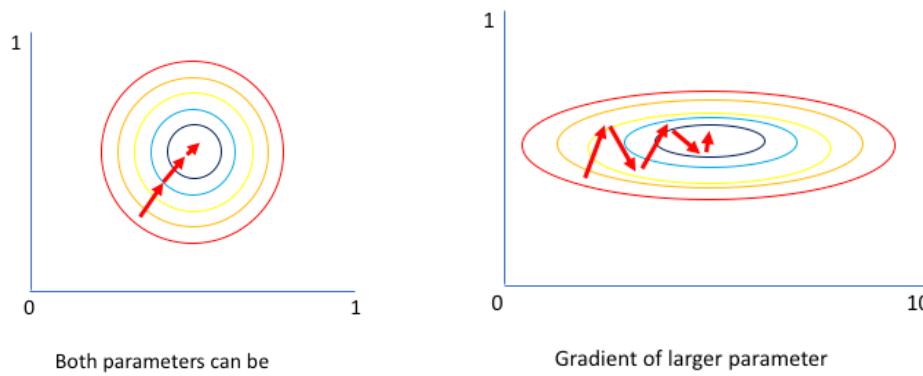
Add

- "Skip connections" aka "residual blocks"
 - $\text{output} = \text{module}(\text{input}) + \text{input}$
 - Allows gradient to flow from the loss function all the way to the first layer
 - (Possible because each module's output is the same shape as its input)



Layer Normalization

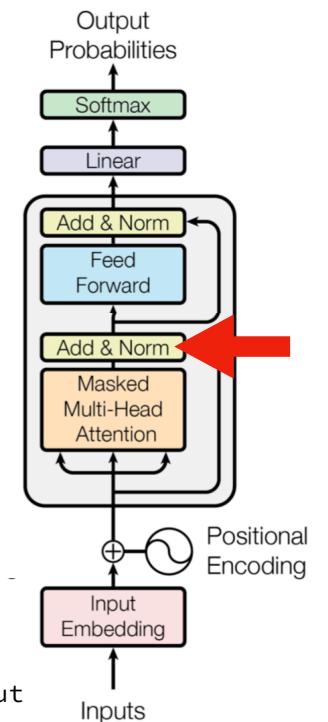
- Neural net modules perform best when input vectors have uniform mean and std in each dimension.
- As inputs flow through the network, means and std's get blown out.
- Layer Normalization is a hack to reset things to where we want them in between layers.



<https://arxiv.org/pdf/1803.08494.pdf>

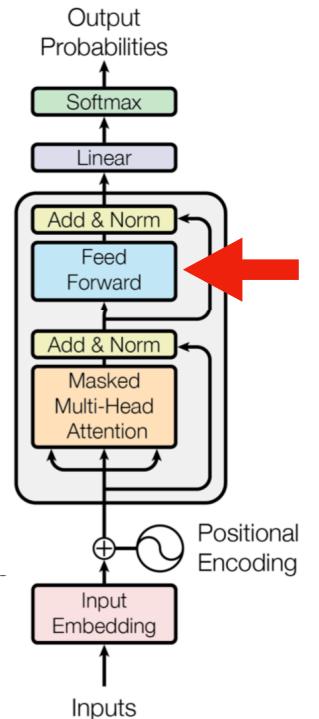
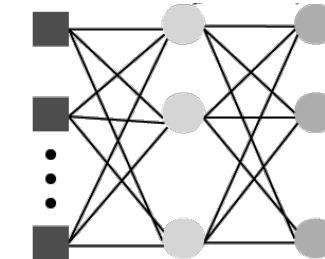
output = module(layer_norm(input)) + input

$$y = \frac{x - E[x]}{\sqrt{Var[x] + \epsilon}}$$



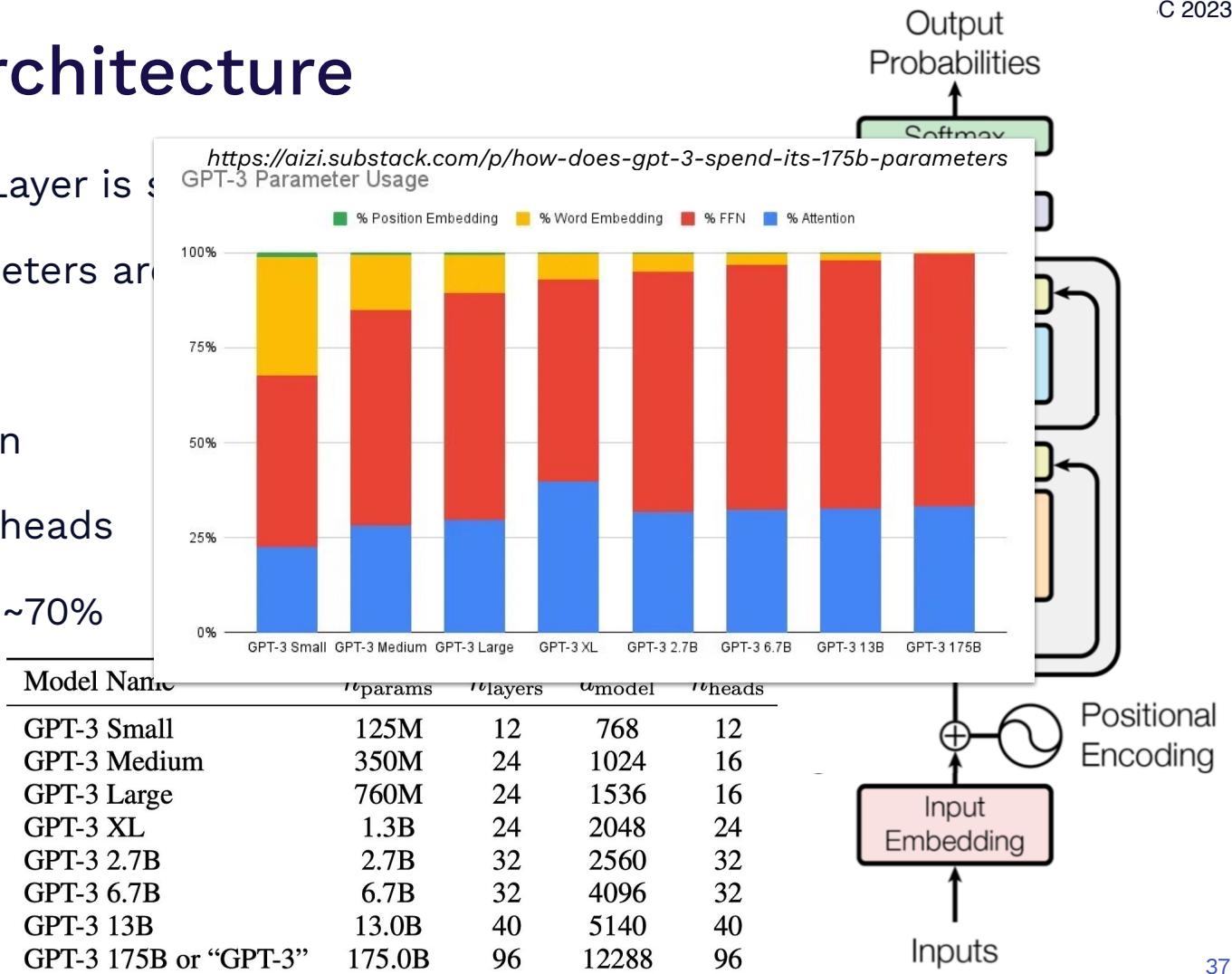
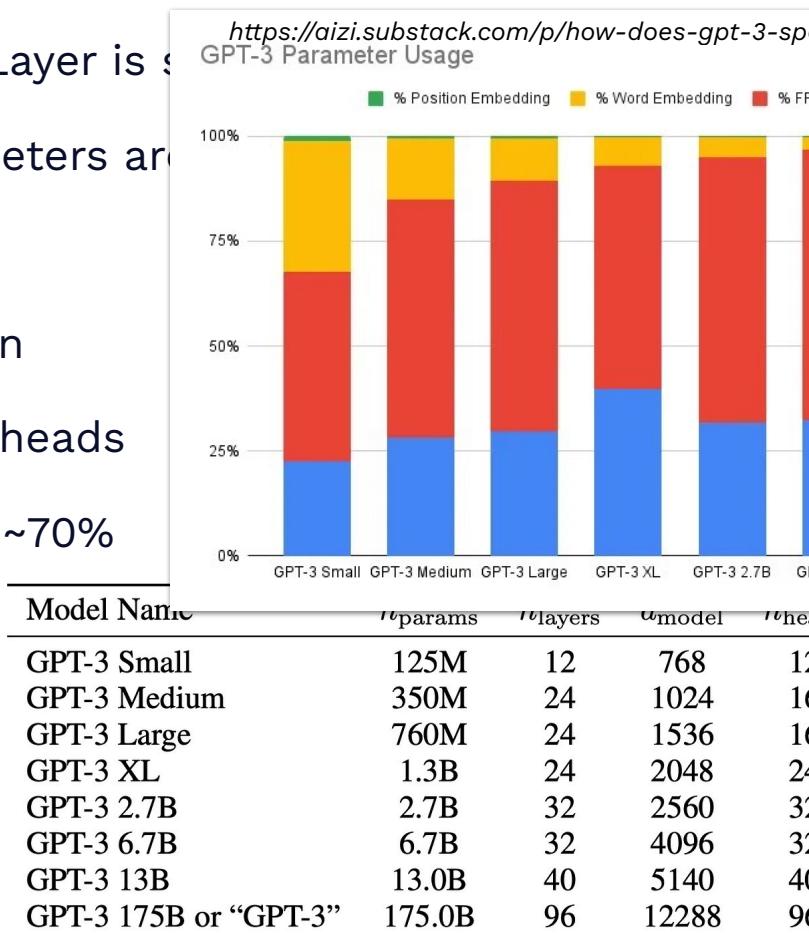
Feed Forward Layer

- Standard Multi-Layer Perceptron with one hidden layer
- Defined $y = W_2(\text{GeLU}(W_1x + b_1)) + b_2$
- Conceptual view:
 - token (augmented with other relevant tokens that it has seen) comes in...
 - ...and "upgrades" its representation



Transformer Architecture

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



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)

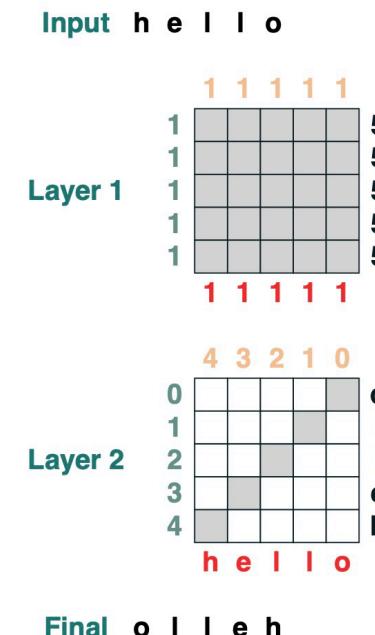
11:54 AM · Oct 19, 2022

490 Retweets 39 Quotes 3,670 Likes 1,023 Bookmarks

Thinking like Transformers

- Restricted Access Sequence Processing (RASP, 2021): programming language of Transformer-implementable operations

```
def flip():
    length = (key(1) == query(1)).value(1)
    flip = (key(length - indices - 1) == query(indices)).value(tokens)
    return flip
flip()
```

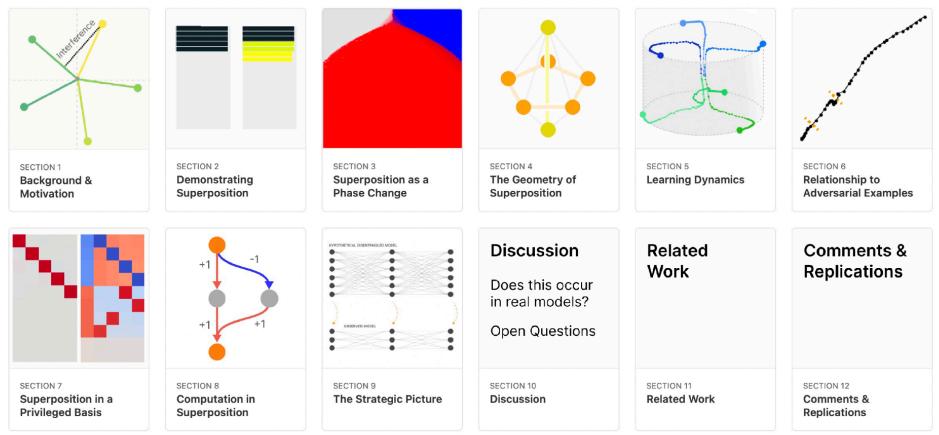


<https://srush.github.io/raspy/>

We mostly don't understand it, though

- Much great work from Anthropic if this has captured your curiosity!

Toy Models of Superposition



AUTHORS

Nelson Elhage*, Tristan Hume*, Catherine Olsson*, Nicholas Schiefer*, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg*, Christopher Olah*

AFFILIATIONS

Anthropic, Harvard

PUBLISHED

Sept 14, 2022

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

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

ONE LAYER
(ATTENTION-ONLY)

Elapsed Training Tokens

0 2.5e9 5.0e9 7.5e9 1e10

-0.1

-0.2

-0.3

-0.4

-0.5

TWO LAYER
(ATTENTION-ONLY)

Elapsed Training Tokens

0 2.5e9 5.0e9 7.5e9 1e10

-0.1

-0.2

-0.3

-0.4

-0.5

THREE LAYER
(ATTENTION-ONLY)

Elapsed Training Tokens

0 2.5e9 5.0e9 7.5e9 1e10

-0.1

-0.2

-0.3

-0.4

-0.5

In-Context Learning Score

We highlight the "phase change" period of training in each plot to make visual comparison easier. The highlighted region is selected for each model based on the derivative of the in-context learning score.

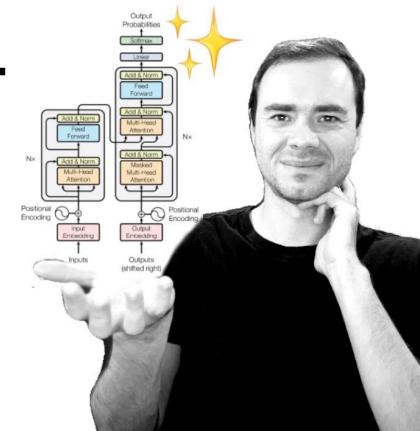
One-layer model
has no sudden improvement.

Models with more than one layer
have a sudden improvement in 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 and MLP blocks

**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?

