



Transformers

Guillem & Roderic, Summer 2025



Transformers

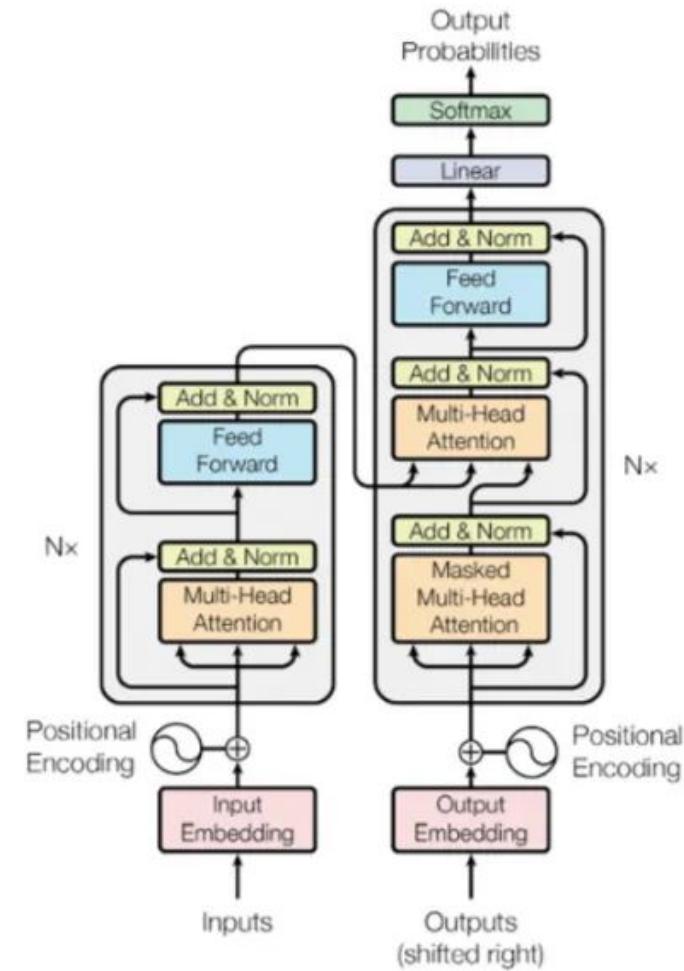
&
Opticons



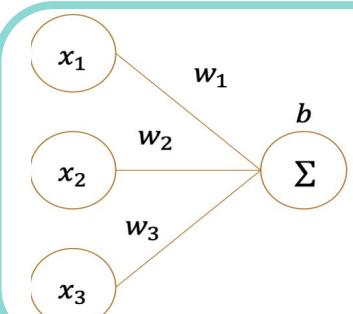


Transformers

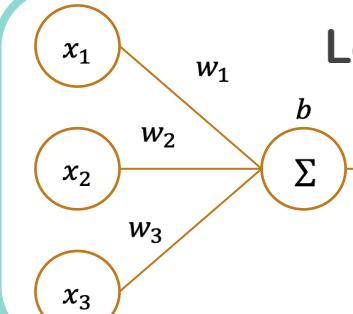
Transformer Attention Is All You Need



Recap:

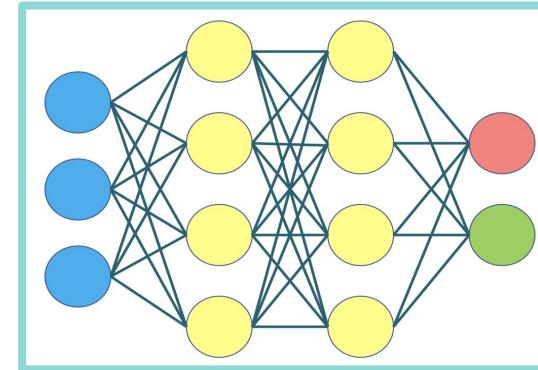


Linear
Regression

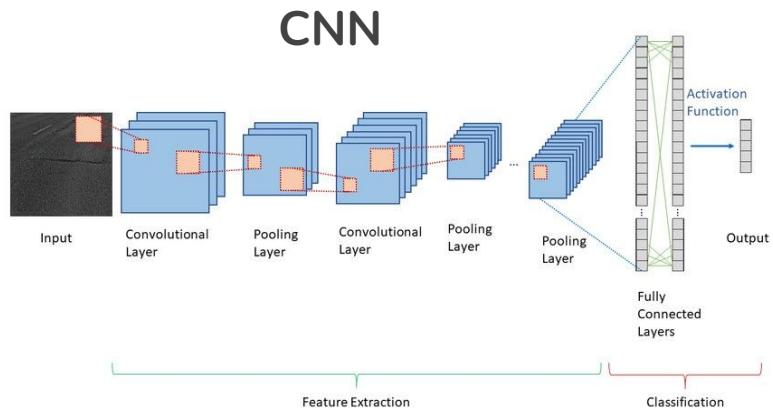


Logistic Regression

Neural Networks

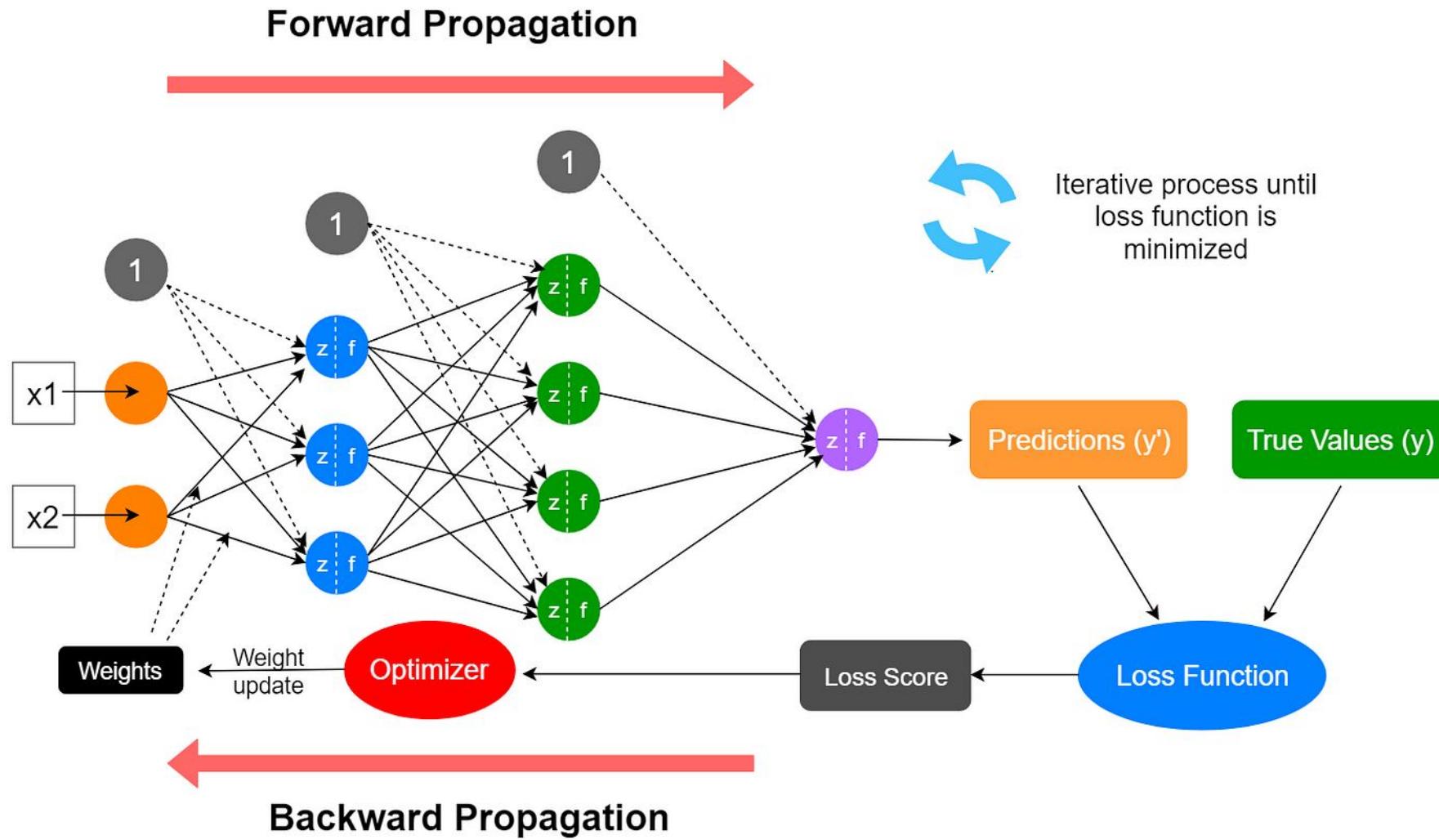


MLP



CNN

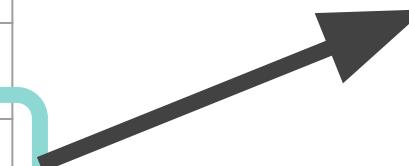
Deep learning training





Vectorization: Numerical Data

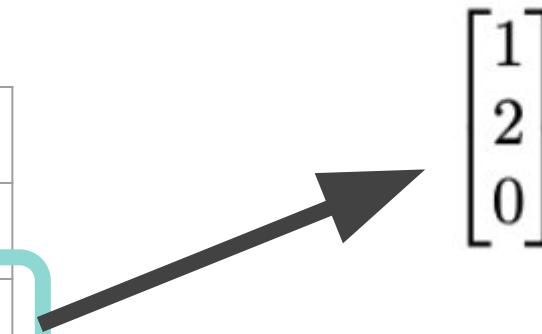
Age	Monthly Salary	Monthly Savings
23	1200	200
34	2000	200
55	1500	300
:	:	:


$$\begin{bmatrix} 34 \\ 2000 \\ 200 \end{bmatrix}$$



Vectorization: Categorical Data

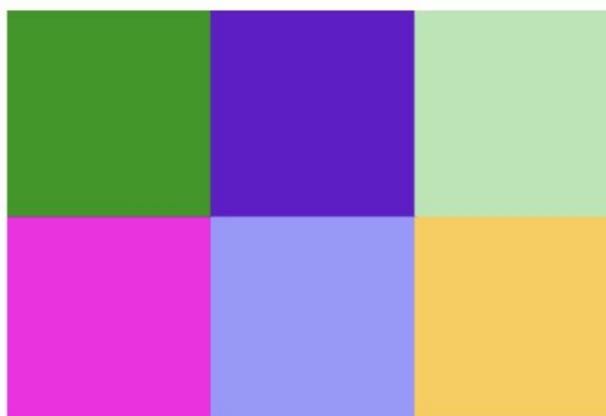
Sex	Hair Color	Smoker
M	Black	Yes
F	Brown	No
F	Black	No
:	:	:



Blond:0, Black:1, Brown: 2, Blue: 3, ...



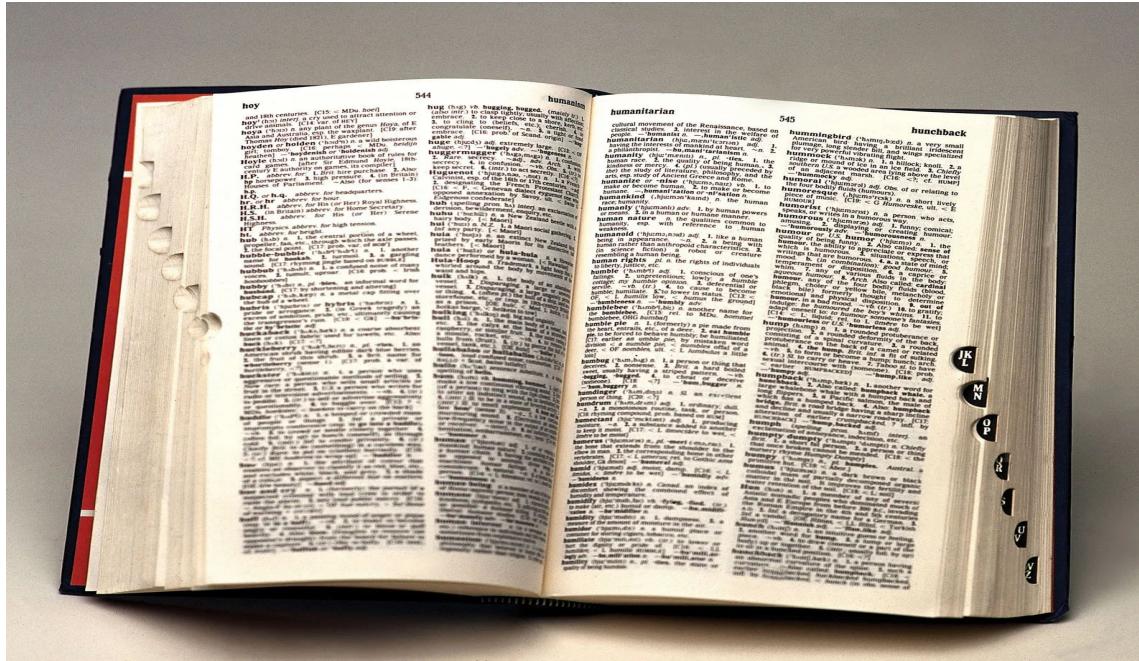
Vectorization: Images



$$\begin{bmatrix} 25 & 51 & 76 \\ 102 & 127 & 153 \\ 178 & 178 & 229 \end{bmatrix}$$

$$\left[\begin{bmatrix} 0 \\ 153 \\ 0 \\ 255 \\ 0 \\ 229 \end{bmatrix}, \begin{bmatrix} 102 \\ 25 \\ 204 \\ 153 \\ 153 \\ 255 \end{bmatrix}, \begin{bmatrix} 178 \\ 229 \\ 178 \\ 255 \\ 204 \\ 76 \end{bmatrix} \right]$$

Vectorization: Words & Text



List of Words in English Dictionary:

A

Aardvark

Abeam

Abacus

Abandon

Abase

...

Vectorization: Words & Text

How can we represent words as vectors?

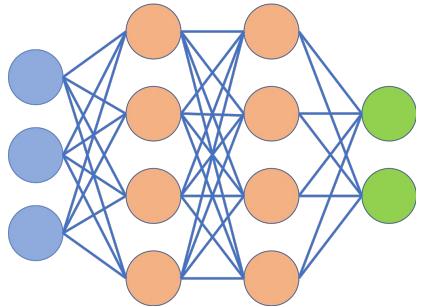
- 1) Assign a number to every word.
- 2) One-hot encoding.
- 3) Words Embeddings

Today!!!

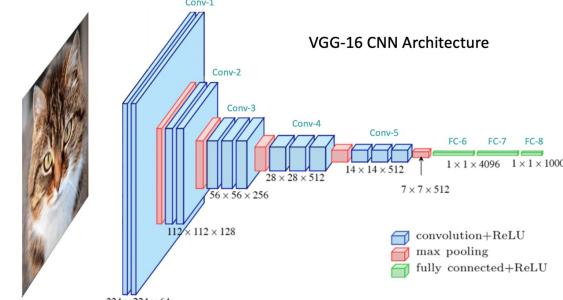
A = 1
Aardvark = 2
Abeam = 3
Abacus = 4
Abandon = 5
Abase = 6
...

A = (1,0,0,0,0,...)
Aardvark = (0,1,0,0,0,...)
Abeam = (0,0,1,0,0,...)
Abacus = (0,0,0,1,0,...)
Abandon = (0,0,0,0,1,...)
Abase = (0,0,0,0,0,1,...)
...

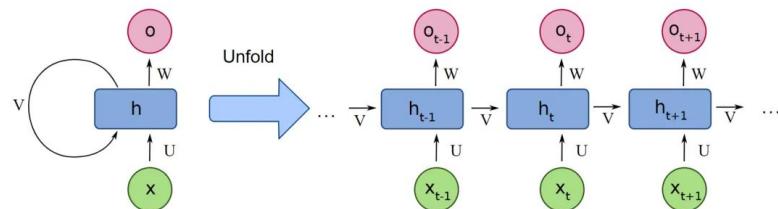
Neural Network architectures



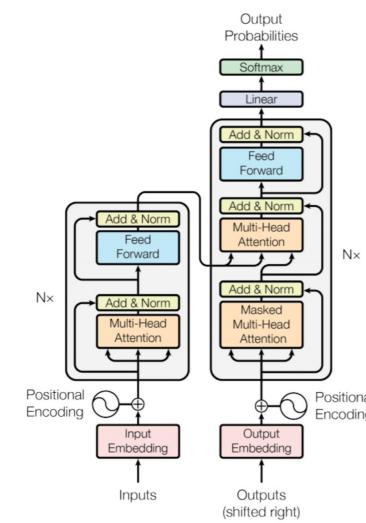
fully connected/linear/dense layers or MLP →
linear map + non-linear activation functions



convolutional layers → CNNs
computer vision, signal processing, etc.



recurrent layers → RNNs, LSTMs, GRUs
time series, natural language processing



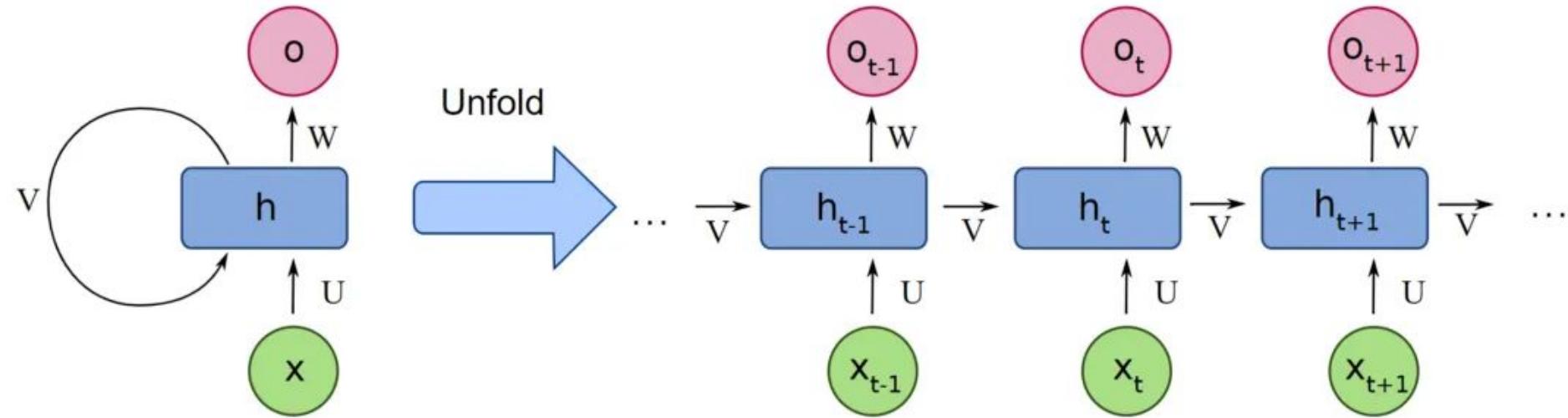
Transformers
Seq to seq models
Multimodal models
Large language models
Generative models
Etc.



NNs in Biology

Type of Data	Example Architecture
Gene expression, single-cell RNA-seq	Fully Connected NN
Biological/medical imaging, histopathology	Convolutional NN
Biological sequences (DNA, RNA, protein), time series	Recurrent Neural Networks, LSTMs
Genomics, proteomics, biomedical language, and structure prediction	Transformers

Recurrent Neural Networks (RNNs)



$x \rightarrow$ feature vector

$h \rightarrow$ hidden State

$U, V, W \rightarrow$ weight matrices

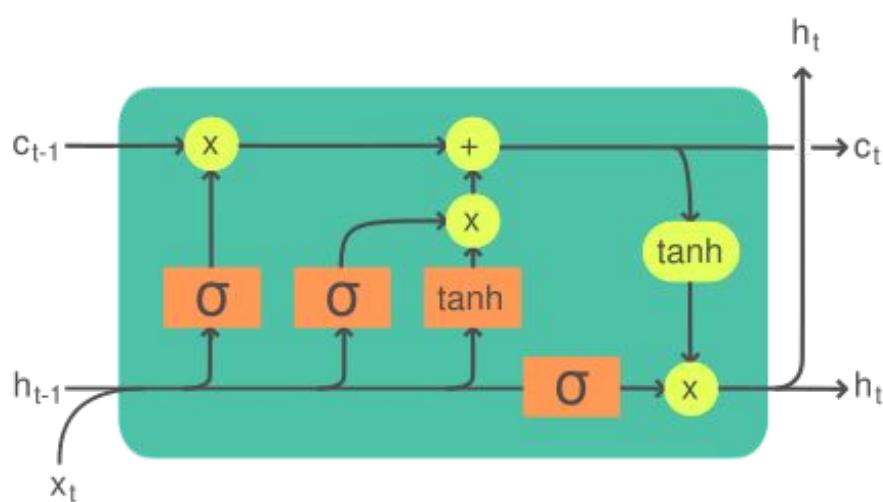
$o \rightarrow$ output

problem: vanishing/exploding gradients

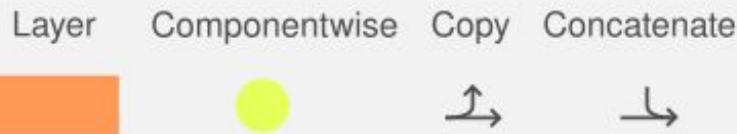
(backpropagation through time)

RNNs lose contextual information rapidly

Long Short-Term Memory (LSTMs)



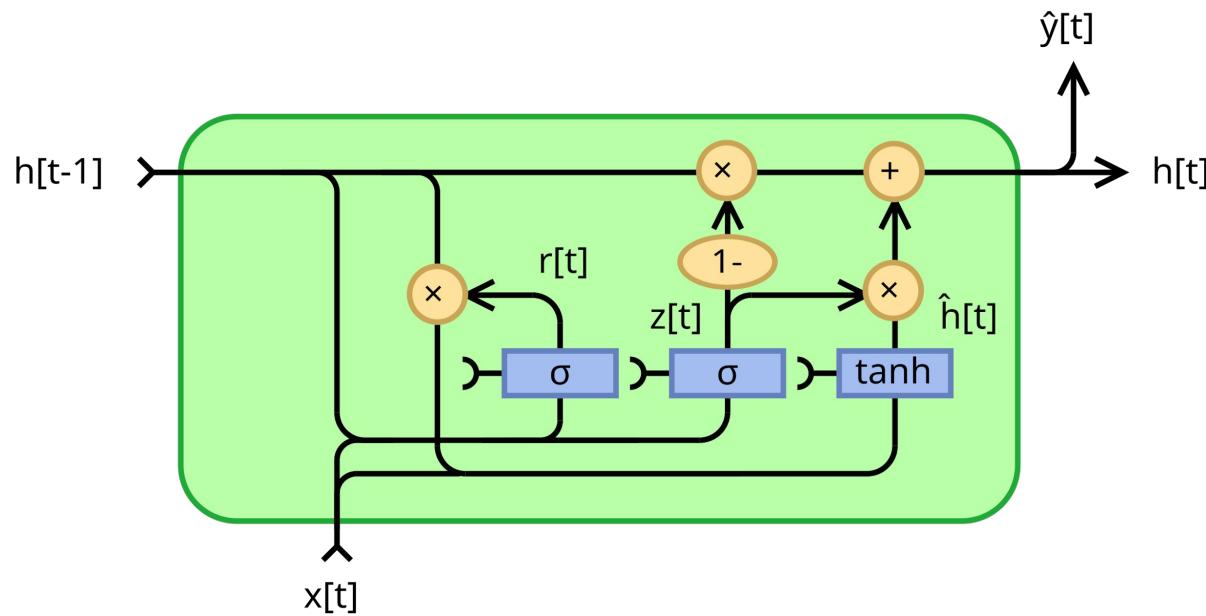
Legend:



$$f_t = \sigma_g(W_f x_t + U_f h_{t-1} + b_f)$$
$$i_t = \sigma_g(W_i x_t + U_i h_{t-1} + b_i)$$
$$o_t = \sigma_g(W_o x_t + U_o h_{t-1} + b_o)$$
$$\tilde{c}_t = \sigma_c(W_c x_t + U_c h_{t-1} + b_c)$$
$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$
$$h_t = o_t \odot \sigma_h(c_t)$$

By adding forget and update gates we can control the information that is retained through time

Gated Recurrent Units (GRUs)



$$z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z)$$

$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r)$$

$$\hat{h}_t = \phi(W_h x_t + U_h (r_t \odot h_{t-1}) + b_h)$$

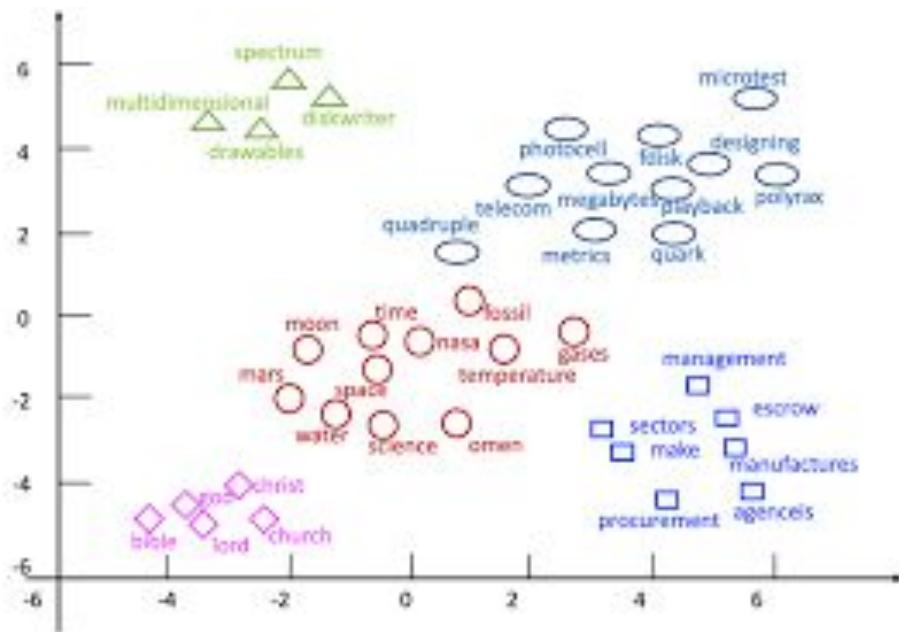
$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \hat{h}_t$$

Like LSTMs but with less parameters

Task: next word prediction

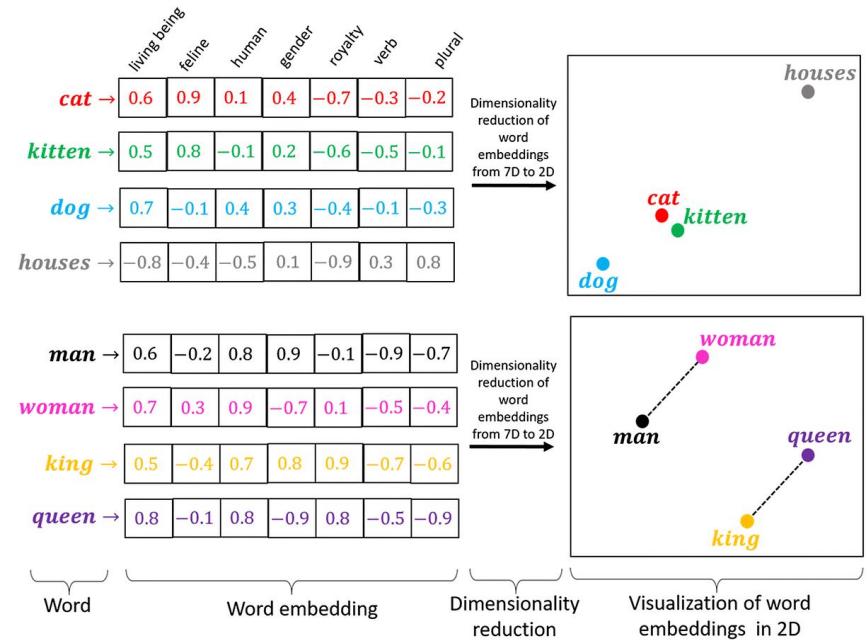
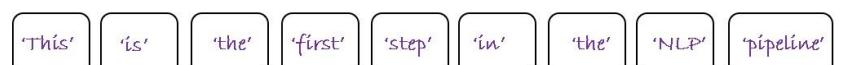
● Tokenization: break text into tokens

- Create word embeddings → words with similar meaning should have similar representations
- Embedding dimension

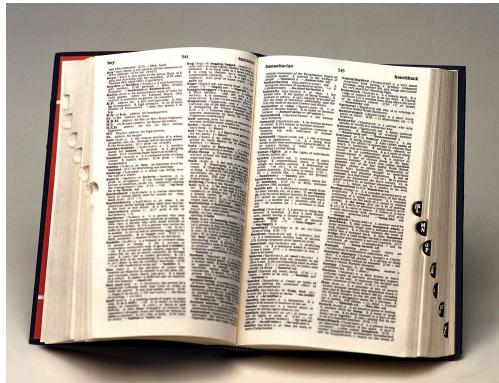


"This is the first step in the NLP pipeline"

Tokenizer



Task: next word prediction



A
Aardvark
Abeam
Abacus
Abandon
...

Embedding dimension

All words, $\sim 50k$

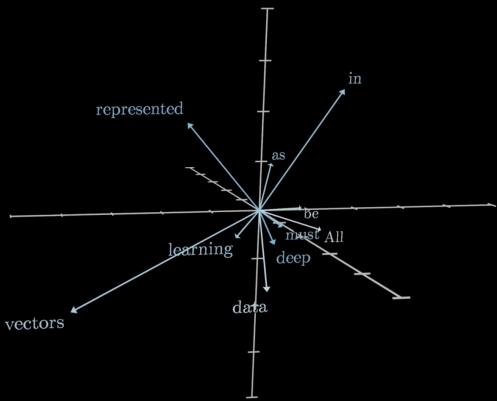
aah	aardvark	aardwolf	aargh	ab	aback	abacterial	abacus	abalone	abandon	...	zygoid	zygomatic	zygomorphic	zygosis	zygote	zygotic	zyme	zymogen	zymosis	zzz
+1.0	+4.3	+2.0	+0.9	-1.5	+2.9	-1.2	+7.8	+9.2	-2.3	...	+0.6	+1.3	+8.4	-8.5	-8.2	-9.5	+6.6	+5.5	+7.3	+9.5
+5.9	-0.8	+5.6	-7.6	+2.8	-7.1	+8.8	+0.4	-1.7	-4.7	...	-0.9	+1.4	-9.5	+2.3	+2.2	+2.3	+8.8	+3.6	-2.8	-1.2
+3.9	-8.7	+3.3	+3.4	-5.7	-7.3	-3.7	-2.7	+1.4	-1.2	...	-7.9	-5.8	-6.7	+3.0	-4.9	-0.7	-5.1	-6.8	-7.7	+3.1
-7.2	-6.0	-2.6	+6.4	-8.0	+6.7	-8.0	+9.4	-0.6	+9.4	...	+4.7	-9.1	-4.3	-7.5	-4.0	-7.5	-3.6	-1.7	-8.6	+3.8
+1.3	-4.6	+0.5	-8.0	+1.5	+8.5	-3.6	+3.3	-7.3	+4.3	...	-6.3	+1.7	-9.5	+6.5	-9.8	+3.5	-4.6	+4.7	+9.2	-5.0
+1.5	+1.8	+1.4	-5.5	+9.0	-1.0	+6.9	+3.9	-4.0	+6.2	...	+7.5	+1.6	+7.6	+3.8	+4.5	+0.0	+9.0	+2.9	-1.5	+2.1
-9.5	-3.9	+3.2	-4.2	+2.3	-1.4	-7.2	-4.0	+1.4	+1.8	...	+3.0	+3.0	-1.4	+7.9	-2.6	-1.3	+7.8	+6.1	+4.0	-7.9
+8.3	+4.2	+9.9	-6.9	+7.3	-6.7	+2.3	-7.4	+6.9	+6.1	...	-1.8	-8.5	+3.9	-0.9	+4.4	+7.3	+9.4	+7.0	-9.7	-2.8
:	:	:	:	:	:	:	:	:	:	...	:	:	:	:	:	:	:	:	:	
-3.7	-2.0	-5.7	-6.2	+8.8	+4.7	-0.2	-5.4	-4.9	-8.8	...	-3.7	+3.9	-2.4	-6.3	-9.4	-8.6	+3.6	-0.9	+0.7	+7.9

Embedding matrix

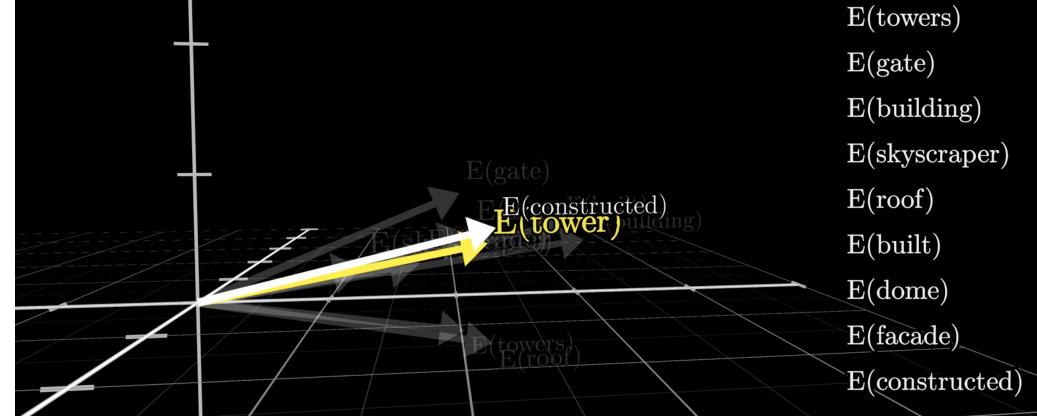


Words $\xrightarrow{\text{"Embedding"}}$ Vectors

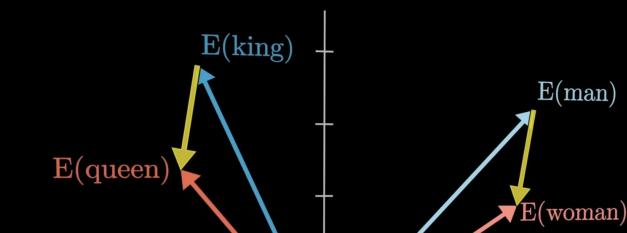
All
data
in
deep
learning
must
be
represented
as
vectors



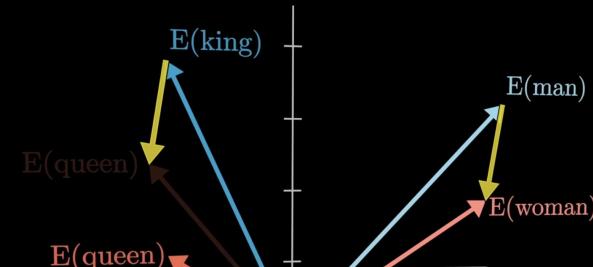
Embeddings closest to E(tower)



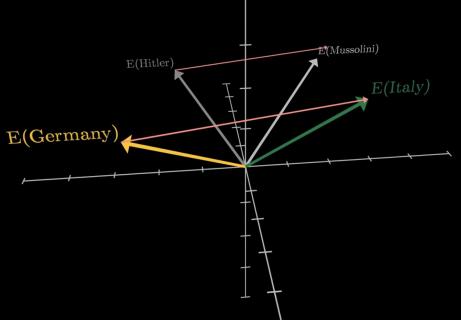
$$E(\text{queen}) \approx E(\text{king}) + E(\text{woman}) - E(\text{man})$$



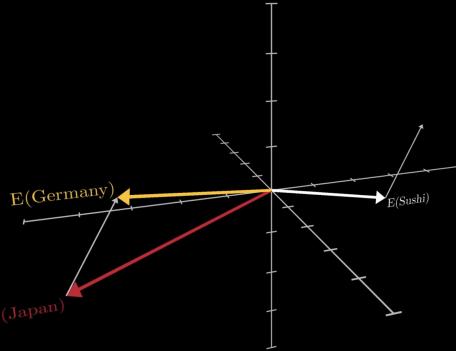
$$E(\text{queen}) \approx E(\text{king}) + E(\text{woman}) - E(\text{man})$$



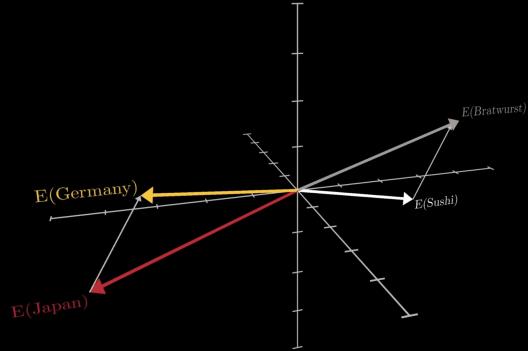
$$E(\text{Hitler}) + E(\text{Italy}) - E(\text{Germany}) \approx E(\text{Mussolini})$$



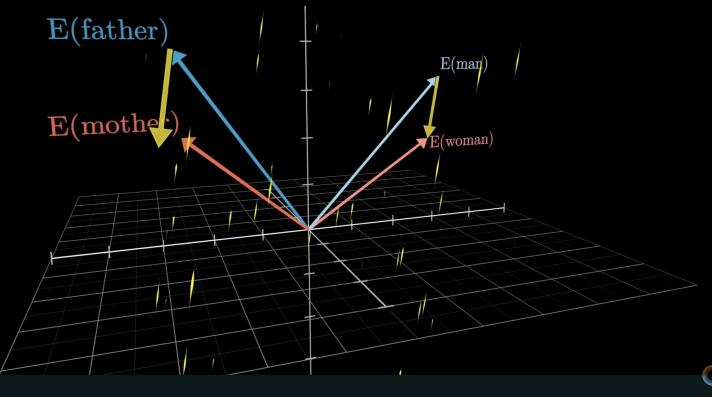
$$[E(\text{Sushi}) + E(\text{Germany}) - E(\text{Japan})] \approx$$



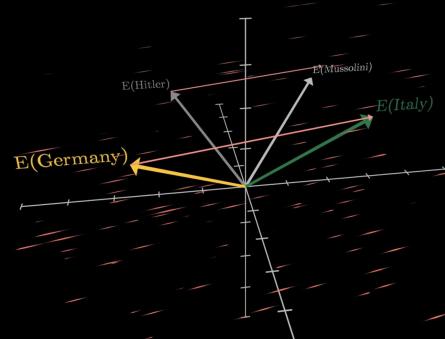
$$E(\text{Sushi}) + E(\text{Germany}) - E(\text{Japan}) \approx [E(\text{Bratwurst})]$$



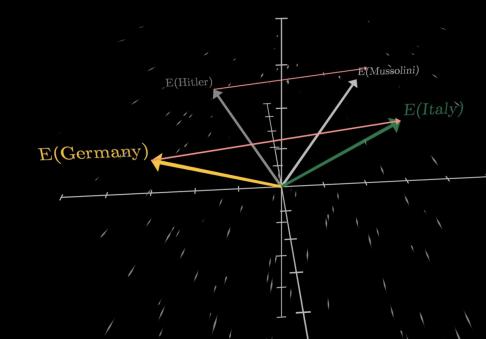
$$E(\text{mother}) - E(\text{father}) \approx E(\text{woman}) - E(\text{man})$$



$$E(\text{Hitler}) + E(\text{Italy}) - E(\text{Germany}) \approx E(\text{Mussolini})$$



$$E(\text{Hitler}) + E(\text{Italy}) - E(\text{Germany}) \approx E(\text{Mussolini})$$



Task: next word prediction

- We can use each word in a sentence to predict the following word; however, we would be missing the context that gives meaning to it
- Context: how the meaning of a word is affected by the meaning of the other words in a given text

1. Machine learning [model](#)
2. Fashion [model](#)

- Language model: from an initial embedding or snippet we want to obtain a single final embedding, which we can then run through an MLP and predict the probability of each word in the vocabulary appearing next

The King doth wake tonight and takes his rouse ...

1 2 3 4 5 6 7 8 9 10

2

King

The King doth wake tonight and takes his rouse ...

King

lived in Scotland

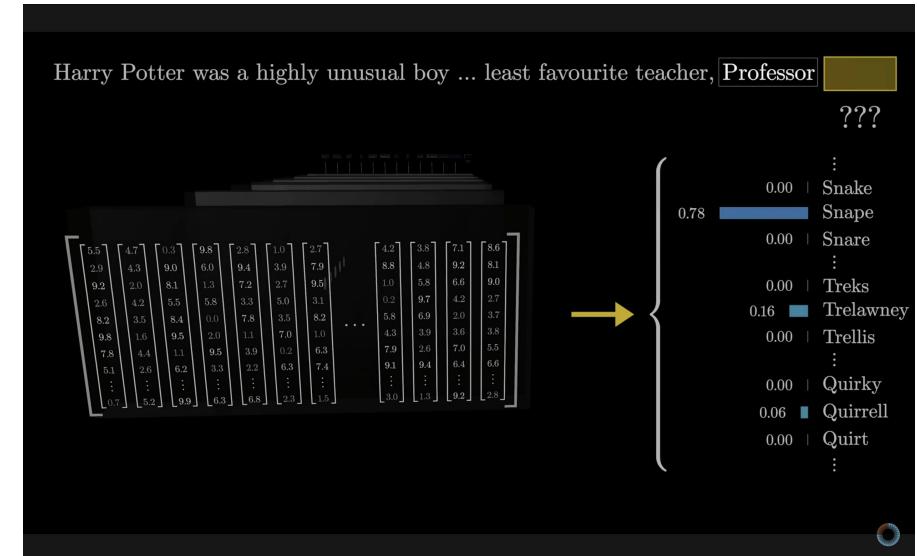
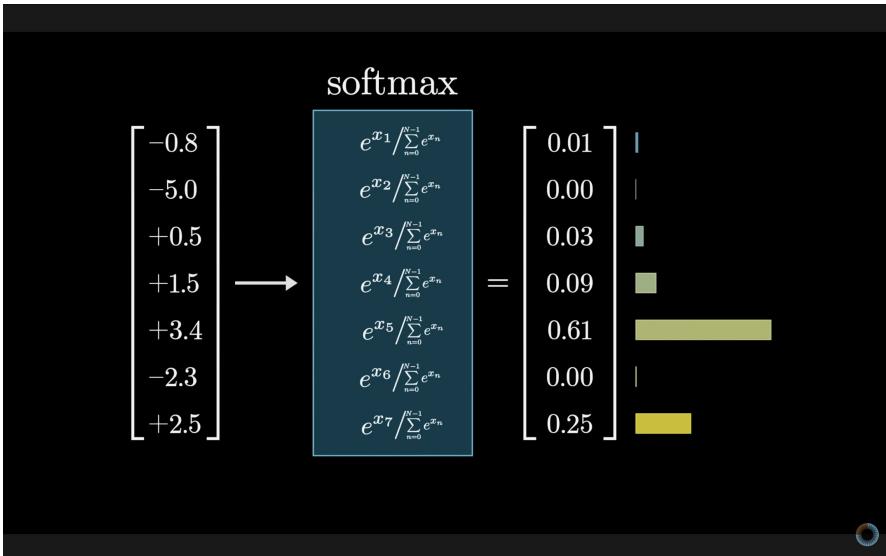
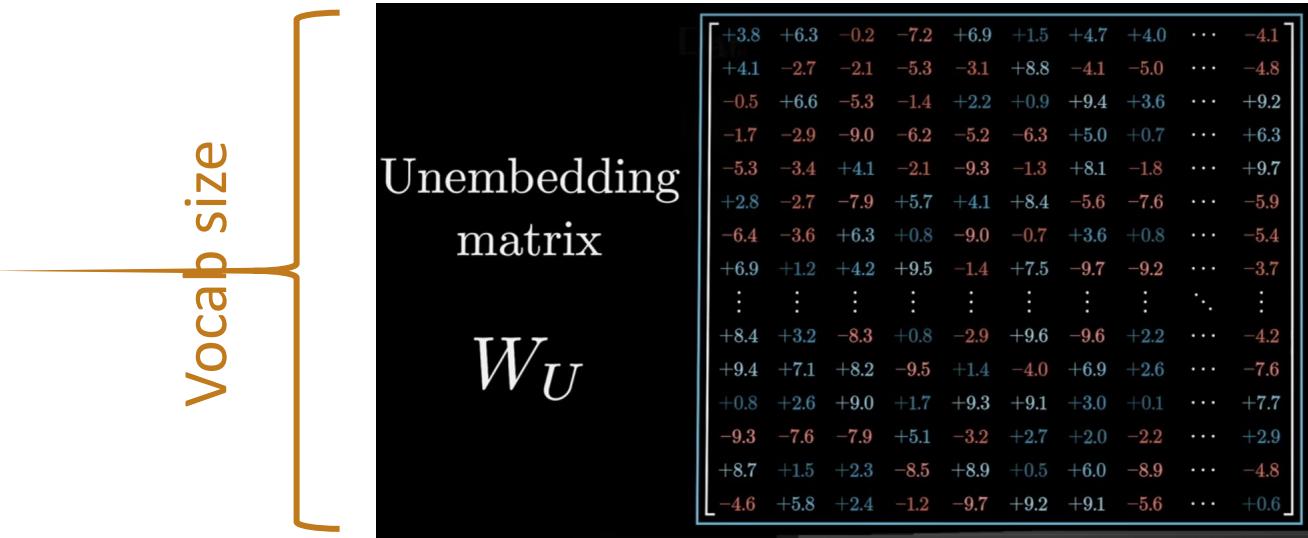
The King doth wake tonight and takes his rouse ...

King
lived in Scotland
murdered predecessor

The King doth wake tonight and takes his rouse ...

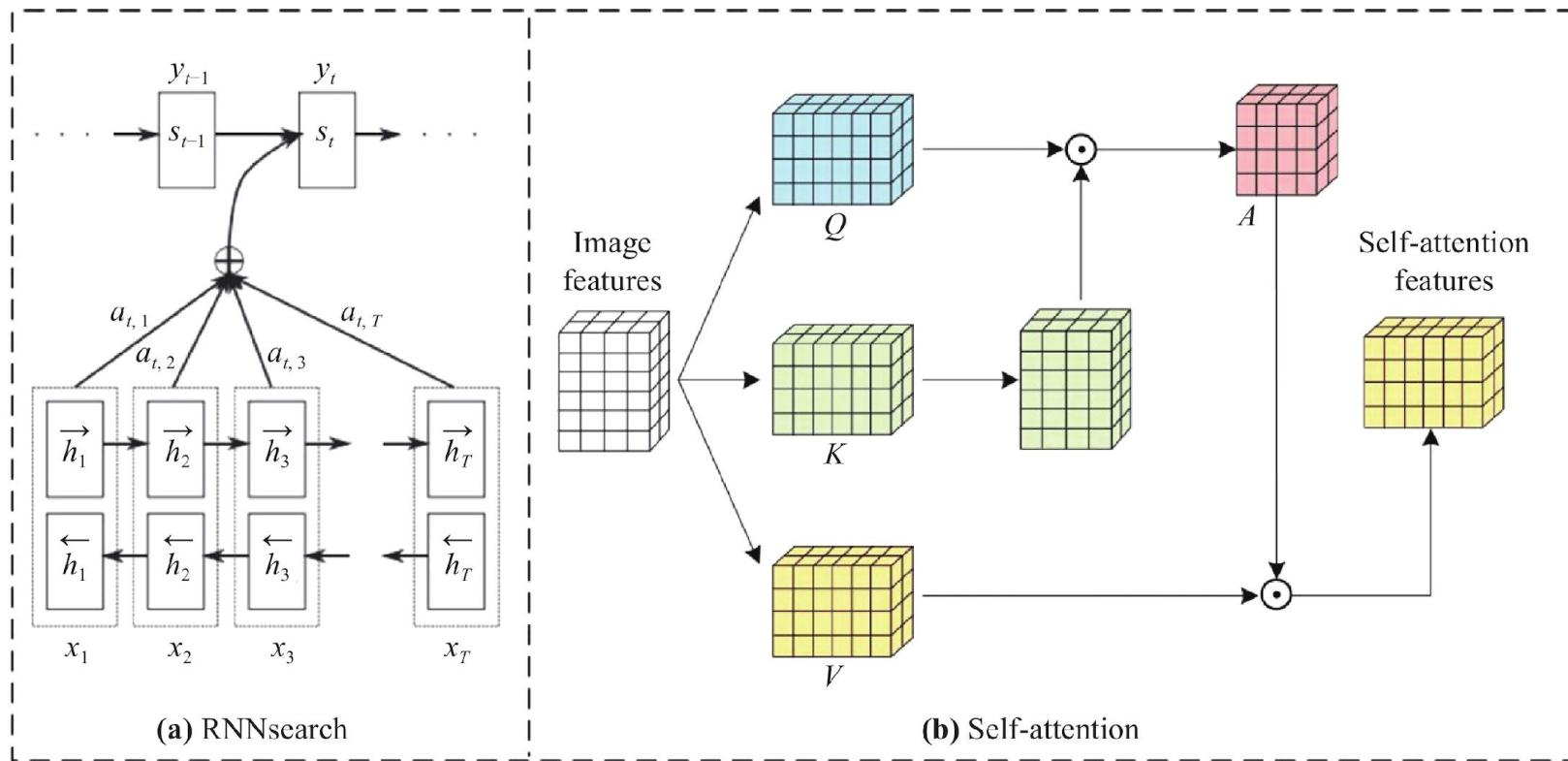
King
lived in Scotland
murdered predecessor
in Shakespearean language

Embedding dimension

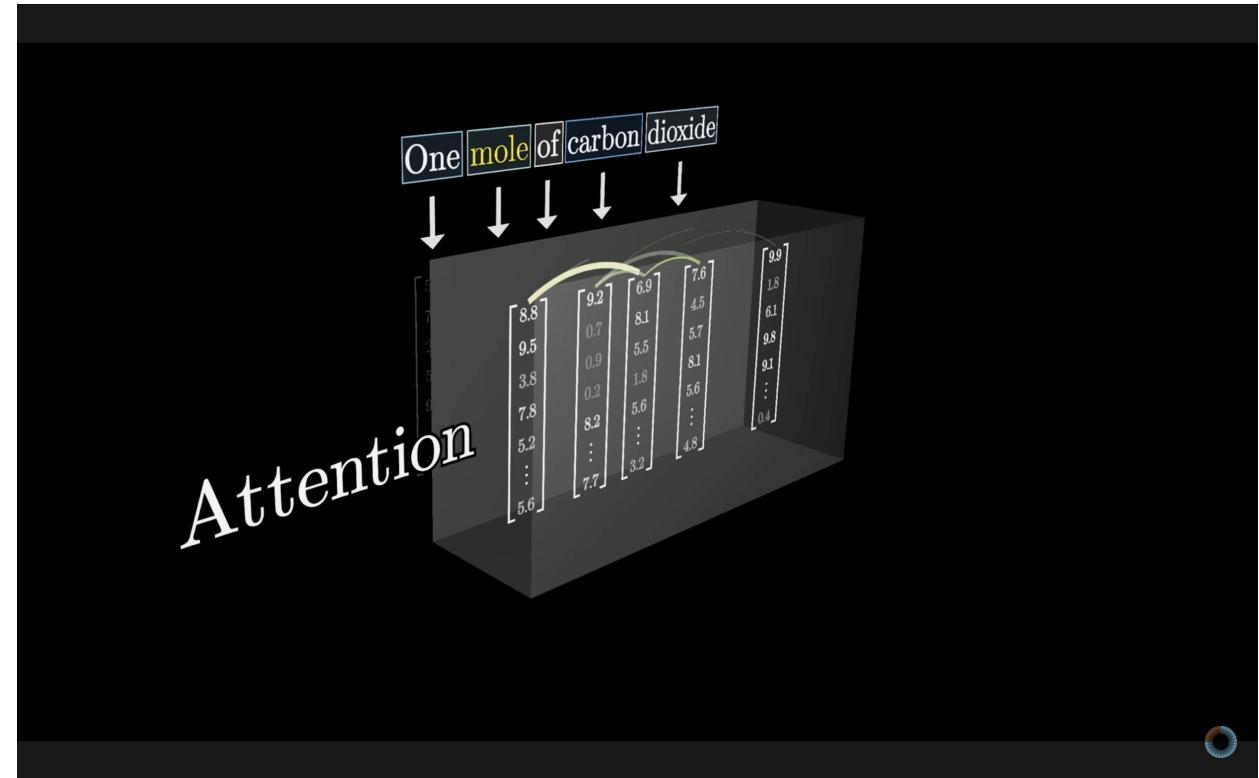
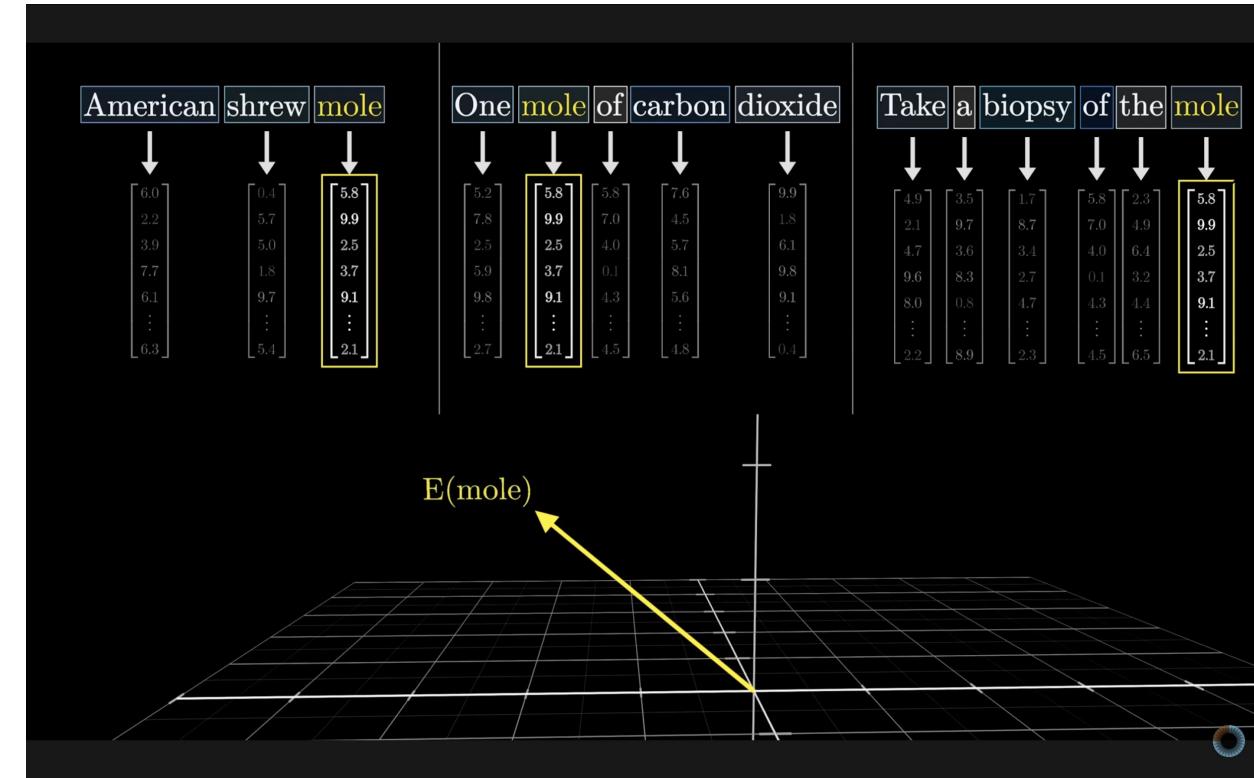


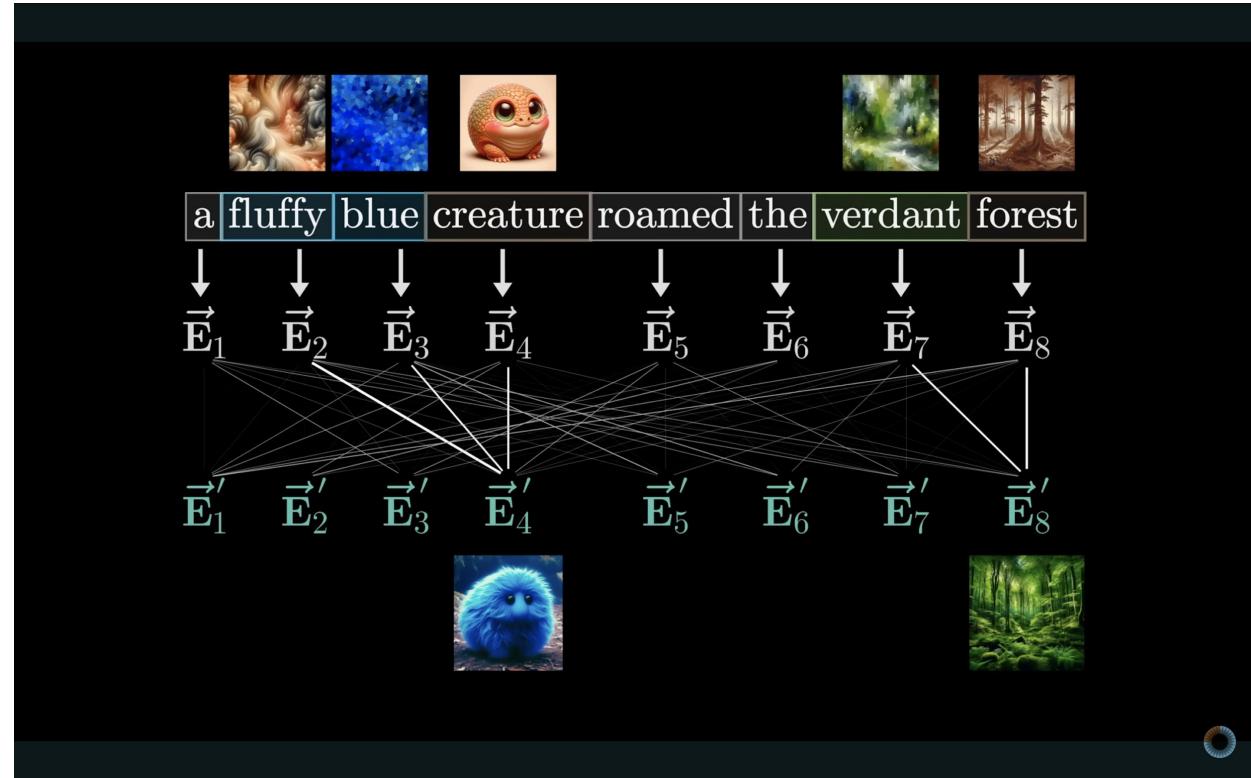
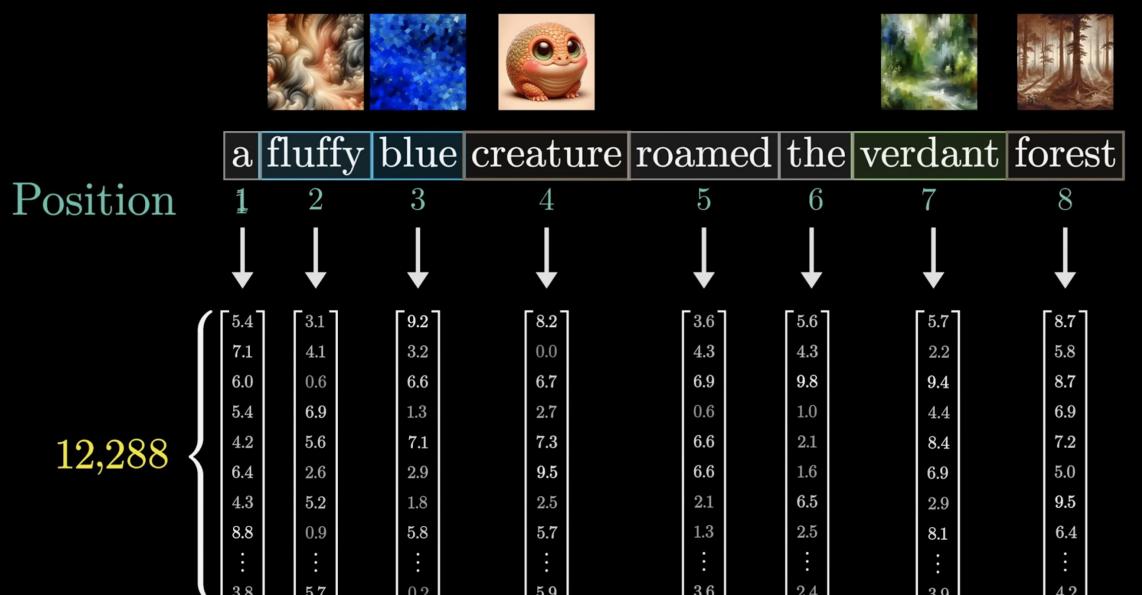
Attention mechanism

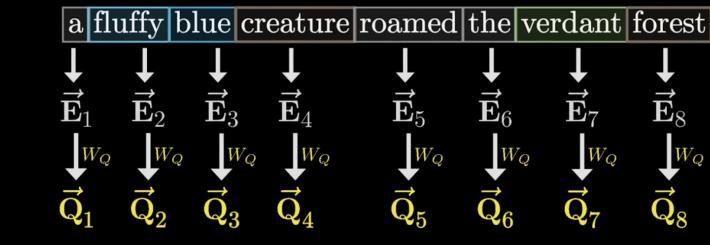
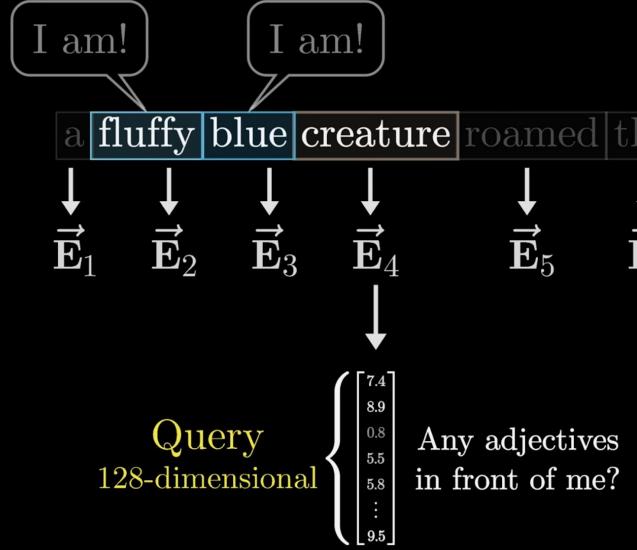
The teacher asked the class a question, and everyone turned to look at her



Li, X., Li, M., Yan, P., Li, G., Jiang, Y., Luo, H., & Yin, S. Deep Learning Attention Mechanism in Medical Image Analysis: Basics and Beyonds. *International Journal of Network Dynamics and Intelligence*. 2023, 2(1), 93–116.

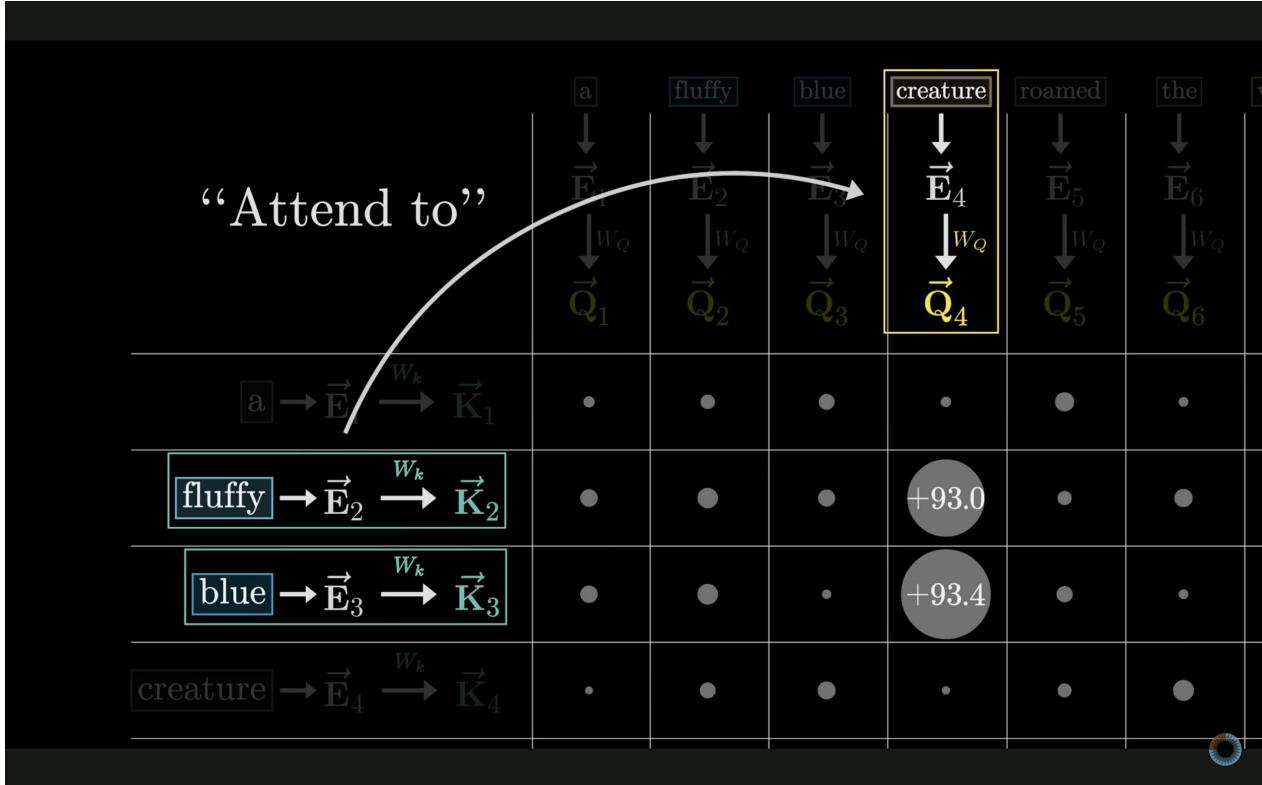


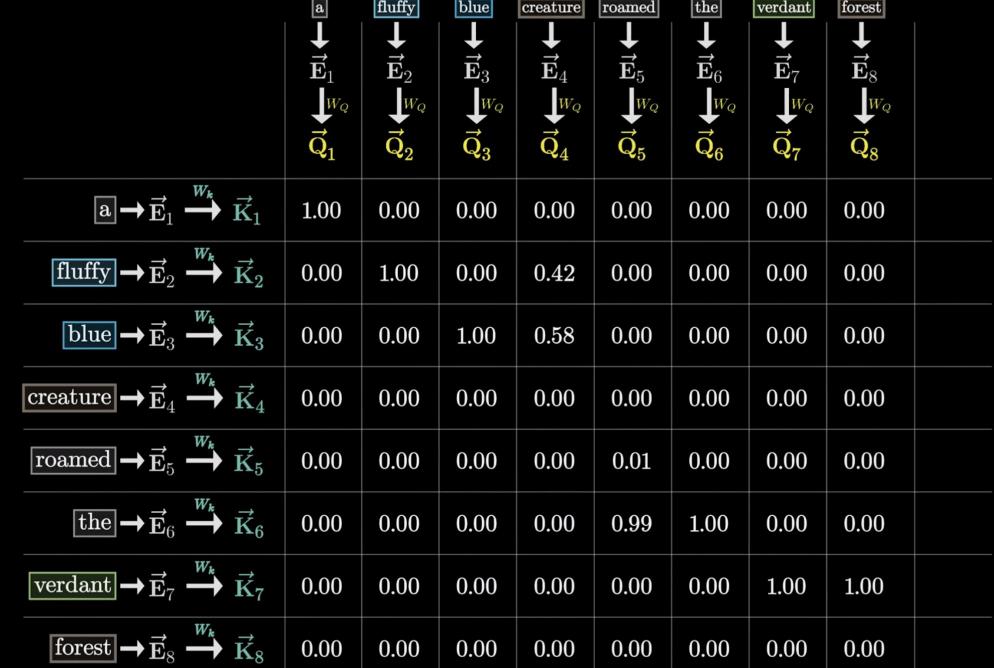
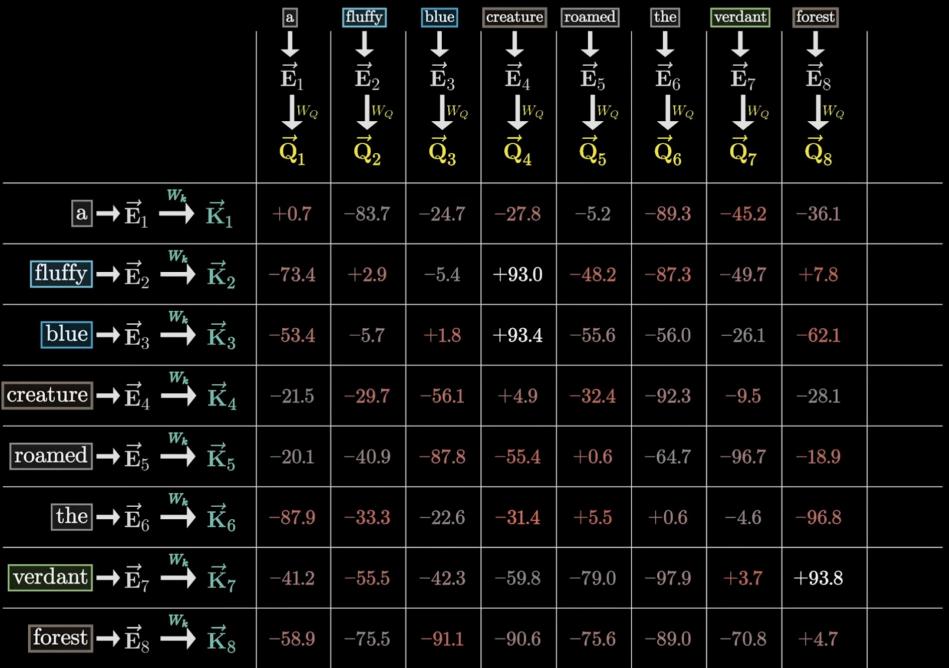




Any adjectives
in front of me?

$$\underbrace{\begin{bmatrix} +0.9 & -0.9 & +7.6 & -0.8 & +4.4 & -2.0 & +8.0 & +3.8 & +4.0 & -3.4 & \cdots & +5.1 \\ +2.7 & -5.1 & -6.7 & +5.9 & +9.1 & -0.8 & +1.8 & +7.1 & -0.8 & +8.9 & \cdots & +1.5 \\ +6.4 & +8.1 & +6.2 & -6.7 & +2.6 & -2.0 & -8.7 & -1.5 & -4.8 & +6.9 & \cdots & -0.2 \\ +9.1 & -2.9 & -2.8 & -9.6 & -6.2 & -2.0 & +8.5 & -7.9 & +8.8 & +7.3 & \cdots & -0.9 \\ -3.4 & -5.3 & +2.3 & -9.2 & -9.6 & -1.4 & -8.6 & -4.9 & -5.5 & -4.9 & \cdots & -7.3 \\ \vdots & \ddots & \vdots \\ -9.7 & -7.6 & +2.3 & +9.4 & +2.7 & -1.8 & -6.7 & +2.7 & -0.2 & +9.7 & \cdots & -8.6 \end{bmatrix}}_{W_Q} = \begin{bmatrix} 2.9 \\ 2.4 \\ 1.0 \\ 0.2 \\ 9.2 \\ 6.6 \\ 7.8 \\ 2.8 \\ 5.8 \\ 0.6 \\ \vdots \\ 9.7 \end{bmatrix} = \begin{bmatrix} \vec{E}_i \\ \vec{Q}_i \end{bmatrix}$$





$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

	Q_1	Q_2	Q_3	Q_4	Q_5	\dots	Q_n	
K_1	$\frac{Q_1 \cdot K_1}{\sqrt{d_k}}$	$\frac{Q_2 \cdot K_1}{\sqrt{d_k}}$	$\frac{Q_3 \cdot K_1}{\sqrt{d_k}}$	$\frac{Q_4 \cdot K_1}{\sqrt{d_k}}$	$\frac{Q_5 \cdot K_1}{\sqrt{d_k}}$	\dots	$\frac{Q_n \cdot K_1}{\sqrt{d_k}}$	
K_2	$\frac{Q_1 \cdot K_2}{\sqrt{d_k}}$	$\frac{Q_2 \cdot K_2}{\sqrt{d_k}}$	$\frac{Q_3 \cdot K_2}{\sqrt{d_k}}$	$\frac{Q_4 \cdot K_2}{\sqrt{d_k}}$	$\frac{Q_5 \cdot K_2}{\sqrt{d_k}}$	\dots	$\frac{Q_n \cdot K_2}{\sqrt{d_k}}$	
K_3	$\frac{Q_1 \cdot K_3}{\sqrt{d_k}}$	$\frac{Q_2 \cdot K_3}{\sqrt{d_k}}$	$\frac{Q_3 \cdot K_3}{\sqrt{d_k}}$	$\frac{Q_4 \cdot K_3}{\sqrt{d_k}}$	$\frac{Q_5 \cdot K_3}{\sqrt{d_k}}$	\dots	$\frac{Q_n \cdot K_3}{\sqrt{d_k}}$	
K_4	$\frac{Q_1 \cdot K_4}{\sqrt{d_k}}$	$\frac{Q_2 \cdot K_4}{\sqrt{d_k}}$	$\frac{Q_3 \cdot K_4}{\sqrt{d_k}}$	$\frac{Q_4 \cdot K_4}{\sqrt{d_k}}$	$\frac{Q_5 \cdot K_4}{\sqrt{d_k}}$	\dots	$\frac{Q_n \cdot K_4}{\sqrt{d_k}}$	
K_5	$\frac{Q_1 \cdot K_5}{\sqrt{d_k}}$	$\frac{Q_2 \cdot K_5}{\sqrt{d_k}}$	$\frac{Q_3 \cdot K_5}{\sqrt{d_k}}$	$\frac{Q_4 \cdot K_5}{\sqrt{d_k}}$	$\frac{Q_5 \cdot K_5}{\sqrt{d_k}}$	\dots	$\frac{Q_n \cdot K_5}{\sqrt{d_k}}$	
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\dots	\vdots	

Unnormalized
Attention Pattern

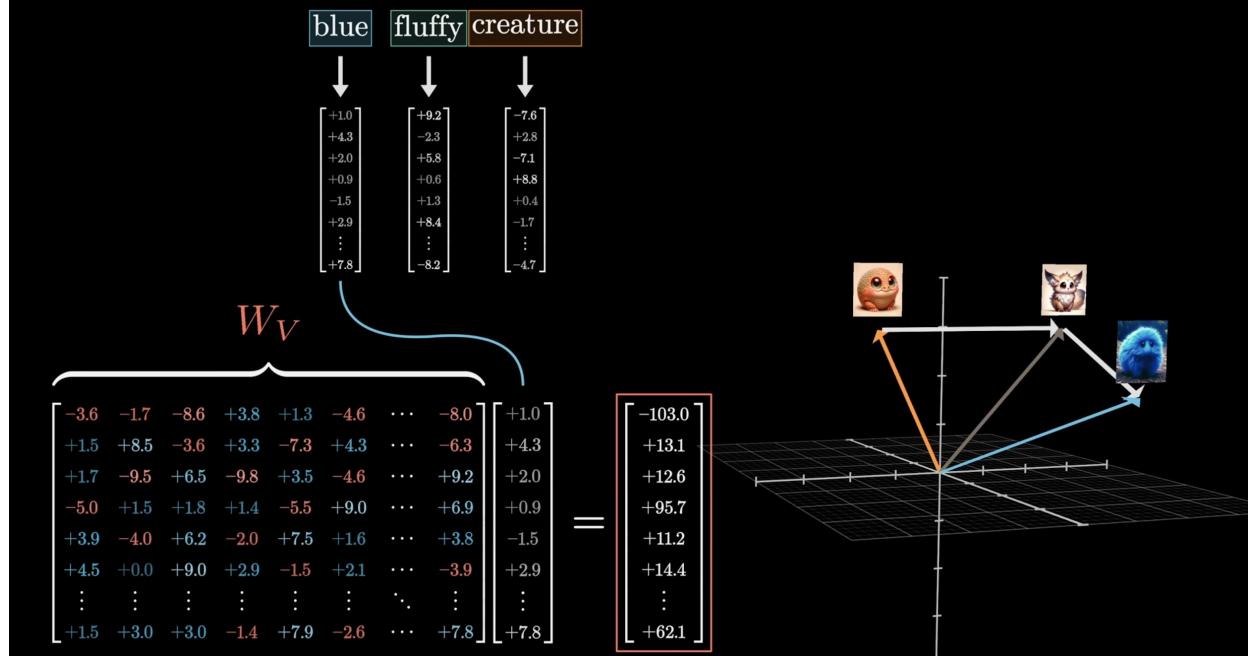
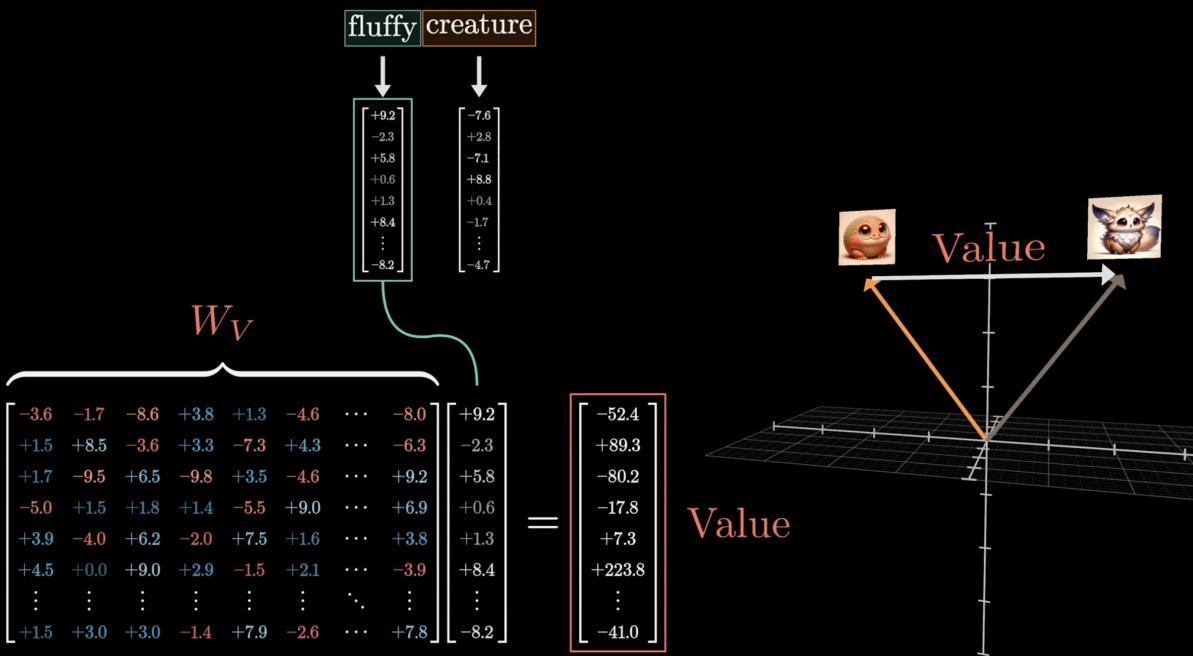
+3.53	+0.80	+1.96	+4.48	+3.74	-1.95
$-\infty$	-0.30	-0.21	+0.82	+0.29	+2.91
$-\infty$	$-\infty$	+0.89	+0.67	+2.99	-0.41
$-\infty$	$-\infty$	$-\infty$	+1.31	+1.73	-1.48
$-\infty$	$-\infty$	$-\infty$	$-\infty$	+3.07	+2.94
$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$	+0.31

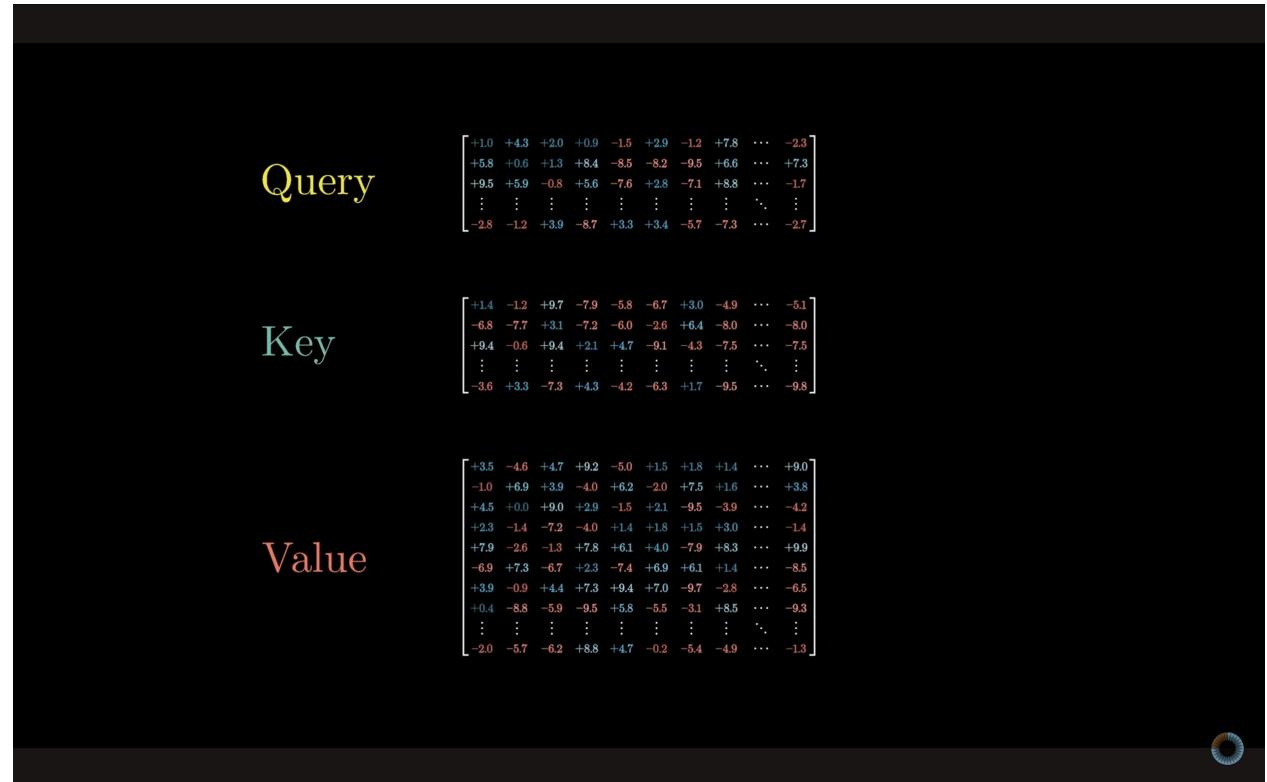
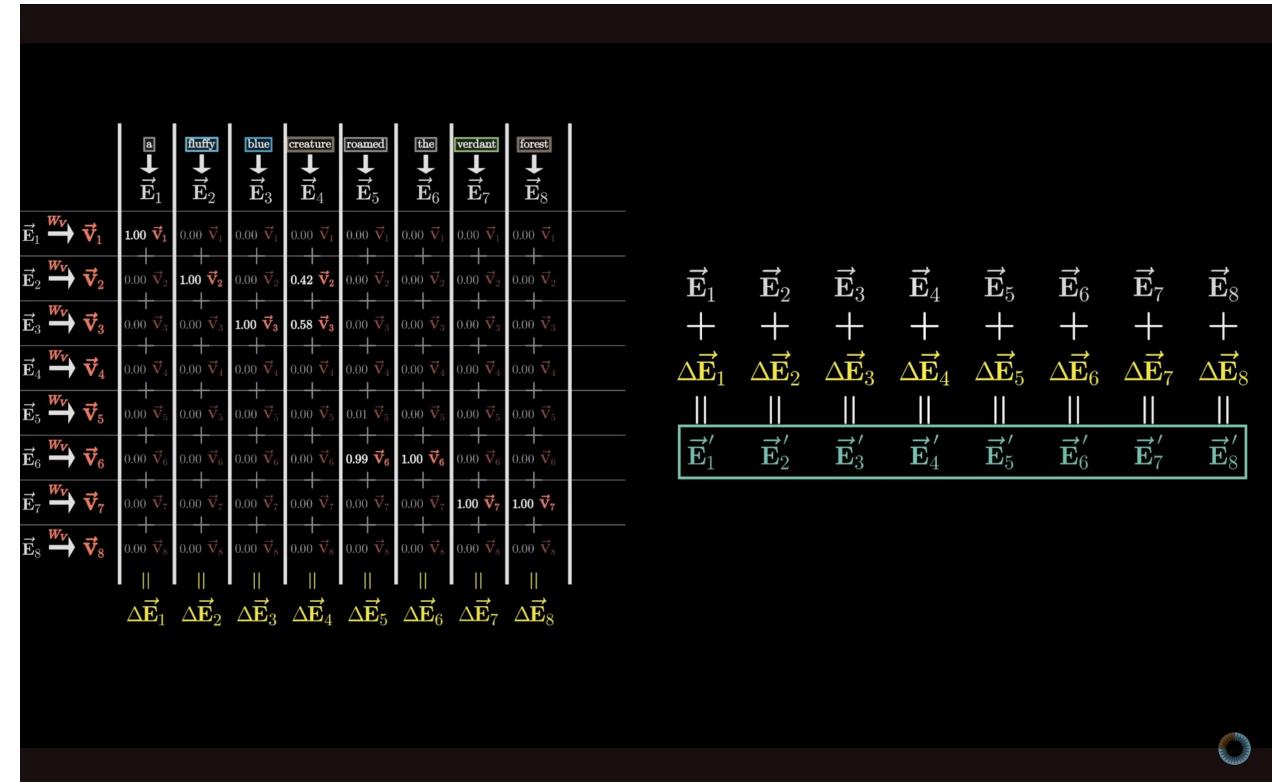
softmax



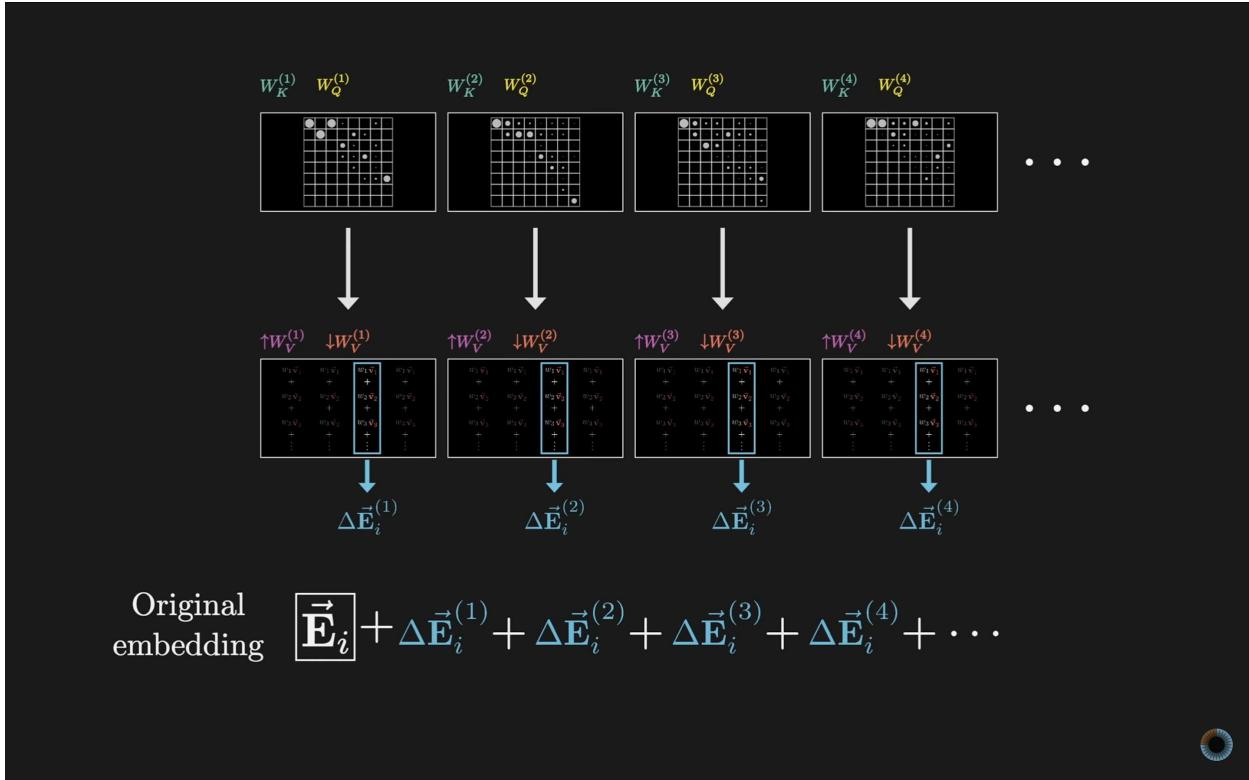
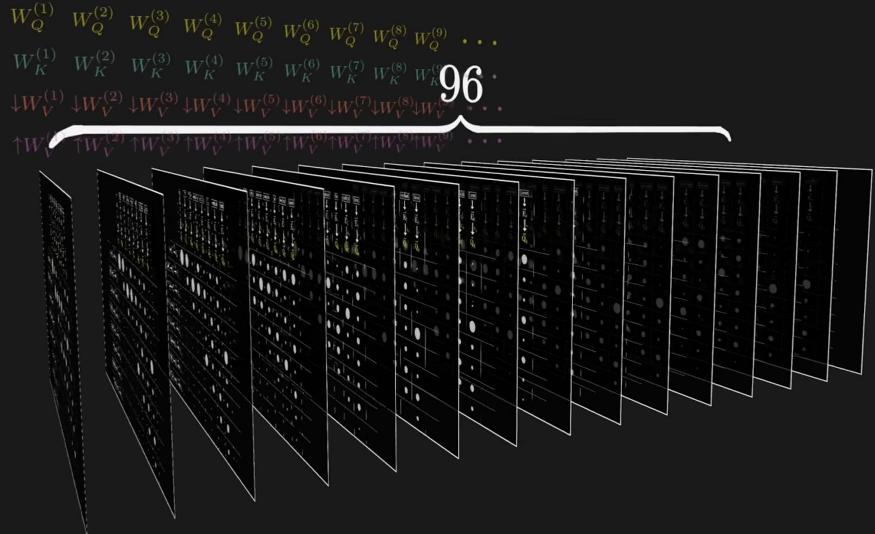
Normalized
Attention Pattern

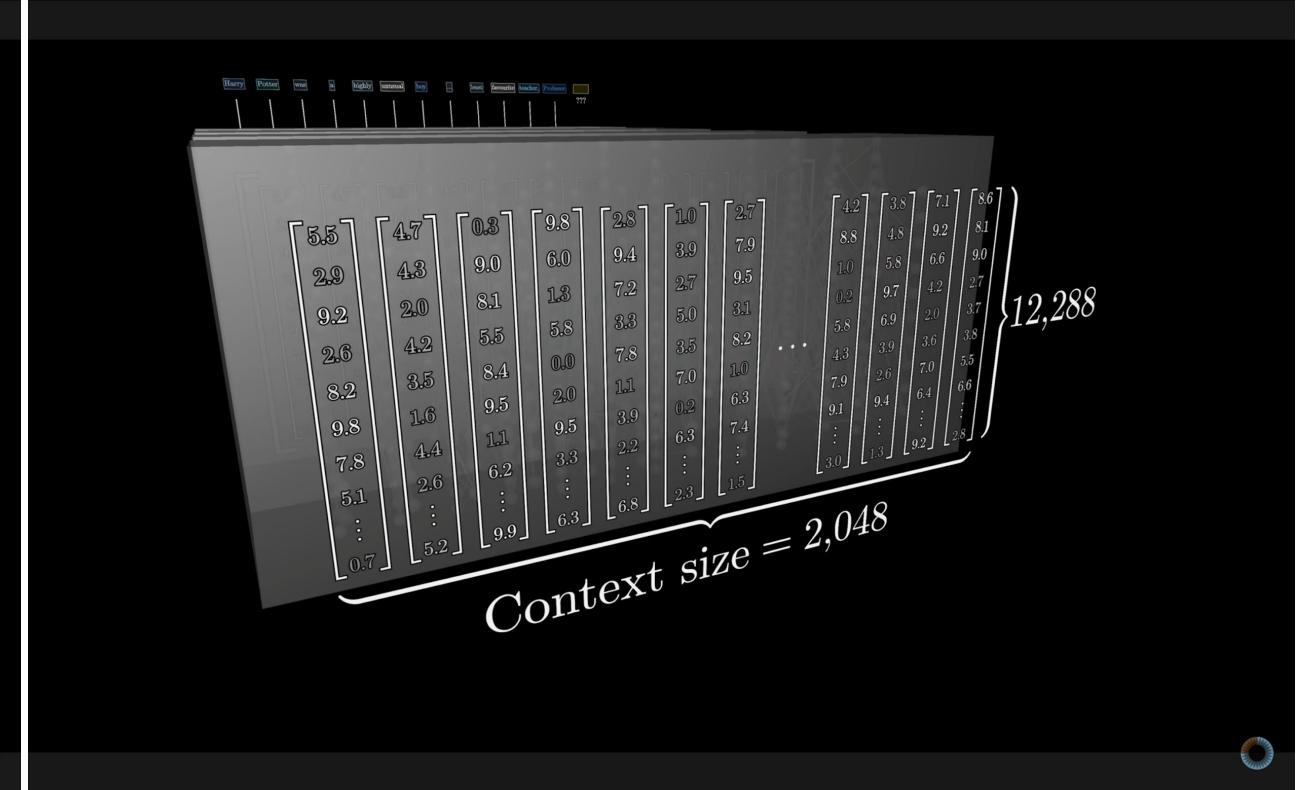
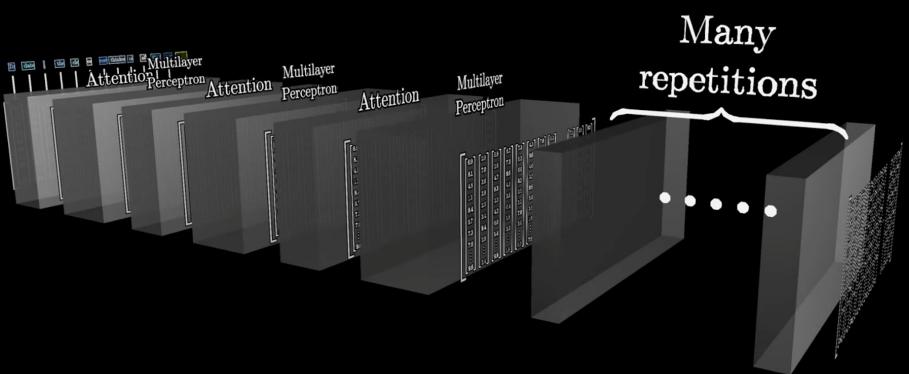
1.00	0.75	0.69	0.92	0.46	0.00
0.00	0.25	0.08	0.02	0.01	0.46
0.00	0.00	0.24	0.02	0.22	0.02
0.00	0.00	0.00	0.04	0.06	0.01
0.00	0.00	0.00	0.00	0.24	0.48
0.00	0.00	0.00	0.00	0.00	0.03





Multi-headed attention





Transformers

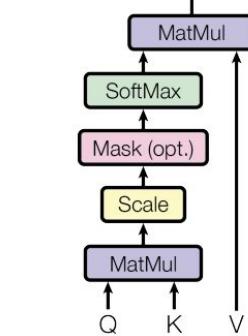
Advantage: we can calculate attention for all the embeddings in a sequence simultaneously. This speeds up computation and allows for massive scalability

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

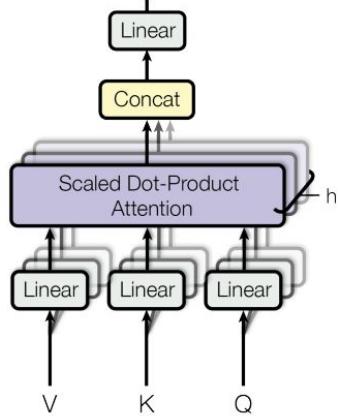
$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

Scaled Dot-Product Attention



Multi-Head Attention



Vaswani, A. "Attention is all you need." *Advances in Neural Information Processing Systems* (2017).

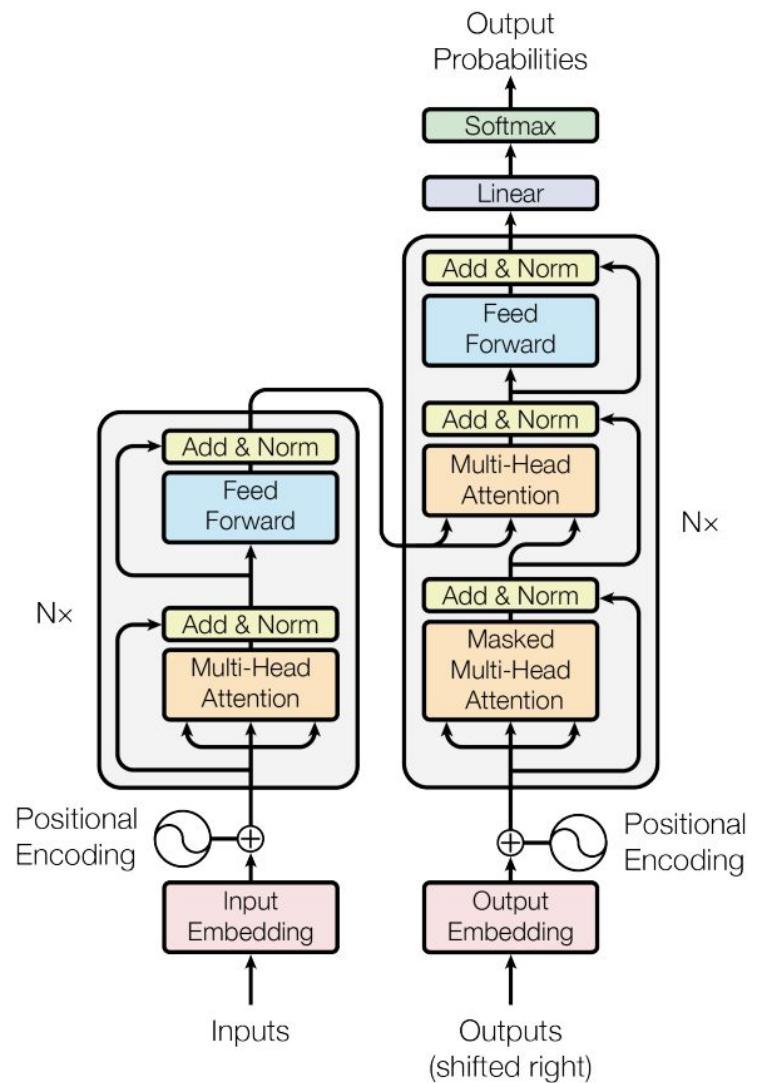
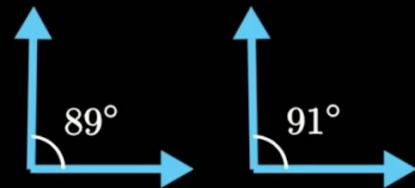


Figure 1: The Transformer - model architecture.

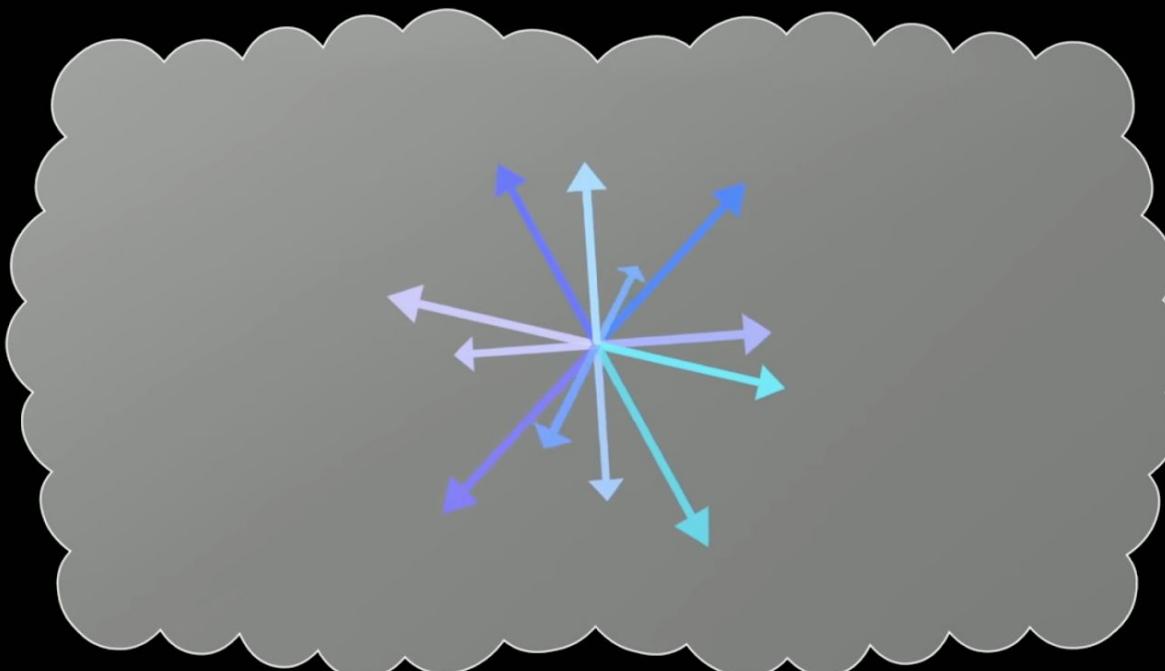


Choose multiple vectors,
each pair between 89° and 91° apart

Johnson–Lindenstrauss
Lemma

\Rightarrow Maximum # of vectors: $\approx \exp(\epsilon \cdot N)$

N -dimensional
Space



YouTube channels



StatQuest with Josh Starmer
@statquest

Neural Networks / Deep Learning

StatQuest with Josh Starmer - 1 / 30



Happy Halloween (Neural Networks Are Not Scary)
StatQuest with Josh Starmer

The Essential Main Ideas of Neural Networks!!!

2 The Essential Main Ideas of Neural Networks

StatQuest with Josh Starmer

The Chain Rule....

3 The Chain Rule

StatQuest with Josh Starmer

Gradient Descent....

4 Gradient Descent, Step-by-Step

StatQuest with Josh Starmer

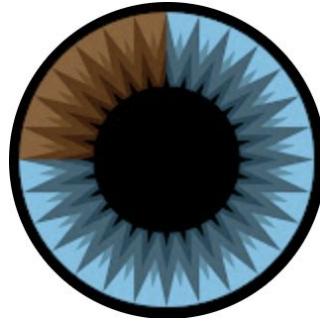
Backpropagation for Neural Networks...

5 Neural Networks Pt. 2: Backpropagation Main Ideas

StatQuest with Josh Starmer

Backpropagation Details...

6 Backpropagation Details Pt. 1: Optimizing 3 parameters...



3Blue1Brown
@3blue1brown

Neural networks

3Blue1Brown - 1 / 8



But what is a neural network? | Deep learning chapter 1
3Blue1Brown

How machines learn

2 How machines learn | DL2

3Blue1Brown

Backpropagation

3 Backpropagation, step-by-step | DL3

3Blue1Brown

Backpropagation calculus

4 Backpropagation calculus | DL4

3Blue1Brown

LLMs

5 Large Language Models explained briefly

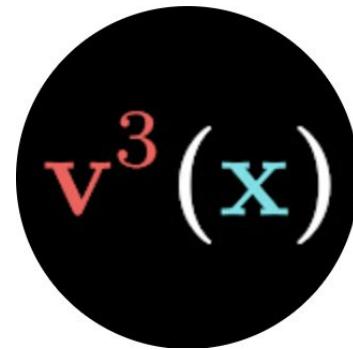
3Blue1Brown

Inside an

6 Transformers (how LLMs work) explained visually | DL5



DeepLearningAI
@Deeplearningai



vcubingx
@vcubingx



Serrano.Academy
@SerranoAcademy