# Transformers

*CS60010*

# Transformers

## Attention Is All You Need (2017)

**Ashish Vaswani**[*]
Google Brain
avaswani@google.com

**Noam Shazeer**[*]
Google Brain
noam@google.com

**Niki Parmar**[*]
Google Research
nikip@google.com

**Jakob Uszkoreit**[*]
Google Research
usz@google.com

**Llion Jones**[*]
Google Research
llion@google.com

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

**Łukasz Kaiser**[*]
Google Brain
lukaszkaiser@google.com

**Illia Polosukhin**[* ‡]
illia.polosukhin@gmail.com

# Great Results with Transformers: Rise of LLMs

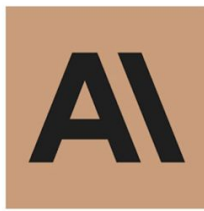Today, Transformer-based models dominate LMSYS Chatbot Arena Leaderboard!

| Rank ▲ | 🏆 Model ▲ | ⭐ Arena Elo ▲ | 📊 95% CI ▲ | 🗳 Votes ▲ | Organization ▲ | License ▲ | Knowledge Cutoff ▲ |
|---|---|---|---|---|---|---|---|
| 1 | GPT-4-Turbo-2024-04-09 | 1258 | +4/-4 | 26444 | OpenAI | Proprietary | 2023/12 |
| 1 | GPT-4-1106-preview | 1253 | +3/-3 | 68353 | OpenAI | Proprietary | 2023/4 |
| 1 | Claude 3 Opus | 1251 | +3/-3 | 71500 | Anthropic | Proprietary | 2023/8 |
| 2 | Gemini 1.5 Pro API-0409-Preview | 1249 | +4/-5 | 22211 | Google | Proprietary | 2023/11 |
| 3 | GPT-4-0125-preview | 1248 | +2/-3 | 58959 | OpenAI | Proprietary | 2023/12 |
| 6 | Meta Llama 3 70b Instruct | 1213 | +4/-6 | 15809 | Meta | Llama 3 Community | 2023/12 |
| 6 | Bard (Gemini Pro) | 1208 | +7/-6 | 12435 | Google | Proprietary | Online |
| 7 | Claude 3 Sonnet | 1201 | +4/-2 | 73414 | Anthropic | Proprietary | 2023/8 |

Gemini / Bard
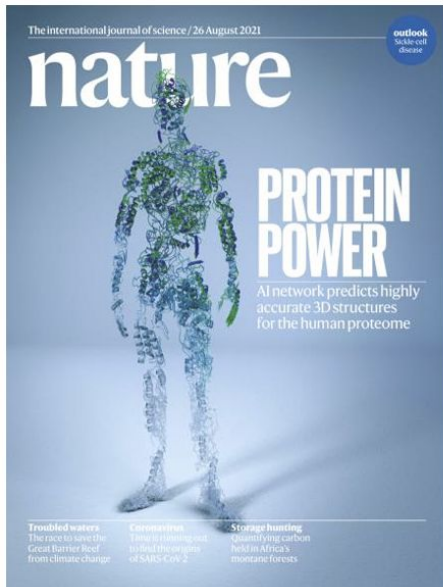(Google)

ChatGPT / GPT-4
(OpenAI)

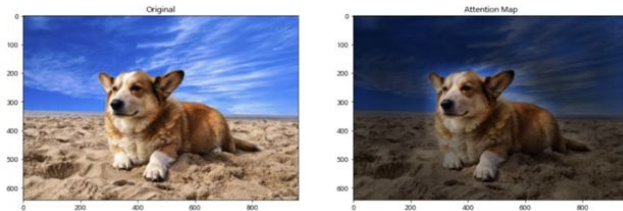Claude 3
(Anthropic)

Llama 3
(Meta)

[Chiang et al., 2024]

https://web.stanford.edu/class/cs224n/

# Transformers have shown promise outside NLP

**Protein Folding**

The international journal of science / 26 August 2021

outlook
Sickle-cell
disease

**nature**

**PROTEIN POWER**

AI network predicts highly
accurate 3D structures
for the human proteome

**Troubled waters**
The race to save the
Great Barrier Reef
from climate change

**Coronavirus**
Threat of reeling out
to find the origins
of SARS-CoV-2

**Storage hunting**
Quantifying carbon
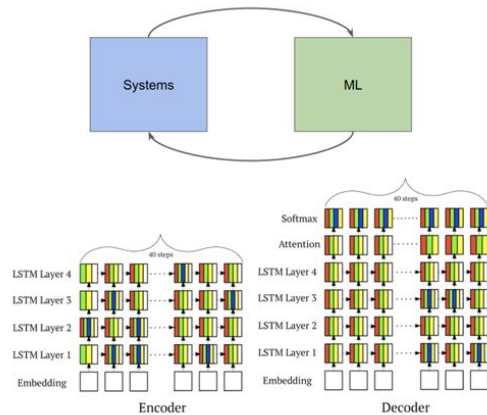held in Africa's
montane forests

[Jumper et al. 2021] aka AlphaFold2!

**Image Classification**

[Dosovitskiy et al. 2020]: Vision Transformer (ViT) outperforms
ResNet-based baselines with substantially less compute.

| | Ours-JFT (ViT-H/14) | Ours-JFT (ViT-L/16) | Ours-I21k (ViT-L/16) | BiT-L (ResNet152x4) | Noisy Student (EfficientNet-L2) |
|---|---|---|---|---|---|
| ImageNet | $88.55 \pm 0.04$ | $87.76 \pm 0.03$ | $85.30 \pm 0.02$ | $87.54 \pm 0.02$ | 88.4/88.5* |
| ImageNet ReaL | $90.72 \pm 0.05$ | $90.54 \pm 0.03$ | $88.62 \pm 0.05$ | 90.54 | 90.55 |
| CIFAR-10 | $99.50 \pm 0.06$ | $99.42 \pm 0.03$ | $99.15 \pm 0.03$ | $99.37 \pm 0.06$ | — |
| CIFAR-100 | $94.55 \pm 0.04$ | $93.90 \pm 0.05$ | $93.25 \pm 0.05$ | $93.51 \pm 0.08$ | — |
| Oxford-IIIT Pets | $97.56 \pm 0.03$ | $97.32 \pm 0.11$ | $94.67 \pm 0.15$ | $96.62 \pm 0.23$ | — |
| Oxford Flowers-102 | $99.68 \pm 0.02$ | $99.74 \pm 0.00$ | $99.61 \pm 0.02$ | $99.63 \pm 0.03$ | — |
| VTAB (19 tasks) | $77.63 \pm 0.23$ | $76.28 \pm 0.46$ | $72.72 \pm 0.21$ | $76.29 \pm 1.70$ | — |
| TPUv3-core-days | 2.5k | 0.68k | 0.23k | 9.9k | 12.3k |

**ML for Systems**

[Zhou et al. 2020]: A Transformer-based
compiler model (GO-one) speeds up a
Transformer model!

| Model (#devices) | GO-one (s) | HP (s) | METIS (s) | HDP (s) | Run time speed up over HP / HDP | Search speed up over HDP |
|---|---|---|---|---|---|---|
| 2-layer RNNLM (2) | 0.173 | 0.192 | 0.355 | 0.191 | 9.9% / 9.4% | 2.95x |
| 4-layer RNNLM (4) | 0.210 | 0.239 | 0.503 | 0.251 | 13.8% / 16.3% | 1.76x |
| 8-layer RNNLM (8) | 0.320 | 0.332 | OOM | 0.764 | 3.8% / 58.1% | 27.8x |
| 2-layer GNMT (2) | 0.301 | 0.384 | 0.344 | 0.327 | 27.6% / 14.3% | 30x |
| 4-layer GNMT (4) | 0.350 | 0.469 | 0.466 | 0.432 | 34% / 23.4% | 58.8x |
| | 0.440 | 0.562 | OOM | 0.693 | 21.7% / 36.5% | 7.35x |
| 2-layer Transformer-XL (2) | 0.223 | 0.268 | 0.37 | 0.262 | 20.1% / 17.4% | 40x |
| 4-layer Transformer-XL (4) | 0.350 | 0.27 | OOM | 0.259 | 17.4% / 12.6% | 26.7x |
| 8-layer Transformer-XL (8) | 0.350 | 0.46 | OOM | 0.425 | 23.9% / 16.7% | 16.7x |
| Inception (2) b64 | 0.229 | 0.312 | OOM | 0.301 | 26.6% / 23.9% | 13.5x |
| AmoebaNet (4) | 0.423 | 0.731 | OOM | 0.498 | 42.1% / 29.3% | 21.0x |
| | 0.394 | 0.44 | 0.426 | 0.418 | 26.1% / 6.1% | 58.8x |
| 2-stack 18-layer WaveNet (2) | 0.317 | 0.376 | OOM | 0.354 | 18.6% / 11.7% | 6.67x |
| 4-stack 36-layer WaveNet (4) | 0.659 | 0.988 | OOM | 0.721 | 50% / 9.4% | 20x |
| GEOMEAN | - | - | - | - | 20.5% / 18.2% | 15x |

Systems

ML

Softmax

Attention

40 steps

LSTM Layer 4

LSTM Layer 3

LSTM Layer 2

LSTM Layer 1

Embedding

Encoder

Decoder

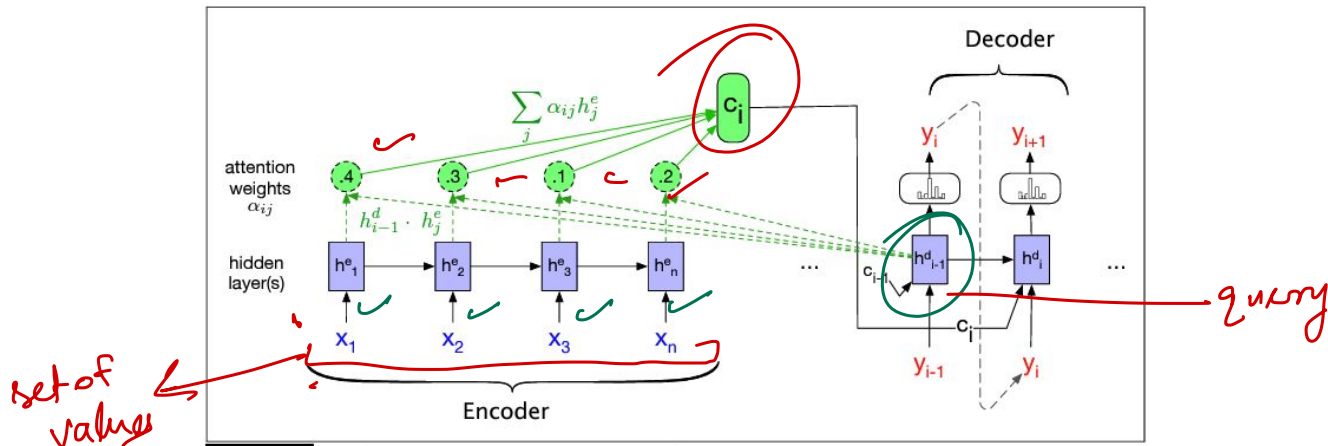# Recap: Attention is a general technique

We can use attention in many architectures not just seq2seq and many tasks (not just MT)

More general definition of attention:

➔ Given a set of vector **values**, and a vector **query**, **attention** is a technique to compute a weighted sum of the values, dependent on the query

We say that the **query attends to the values**

In the seq2seq+attention model, each decoder hidden state (query) *attends to* all the encoder hidden state (values)

# Attention is a general deep learning technique

More general definition of attention:

➔ Given a set of vector **values**, and a vector **query**, <u>**attention**</u> is a technique to compute a weighted sum of the values, dependent on the query
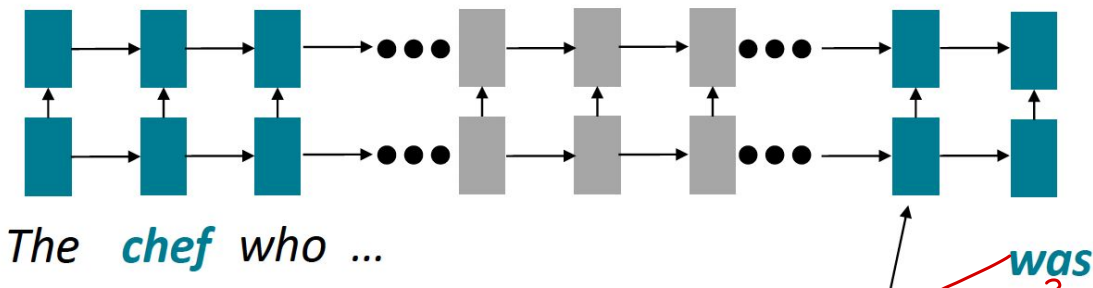
**Intuition:**

- The weighted sum is a **selective summary** of the information contain in the values, where the query determines which values to focus on

- Attention is a way to obtain a **fixed-size representation of an arbitrary set of representations**, *dependent on some other representation*

# Issues with recurrent models

*O(sequence length) steps for distant word pairs*

- Hard to learn long-distance dependencies (gradient problems!)
- Linear order of words is "baked in"; not the right way to think about sentences
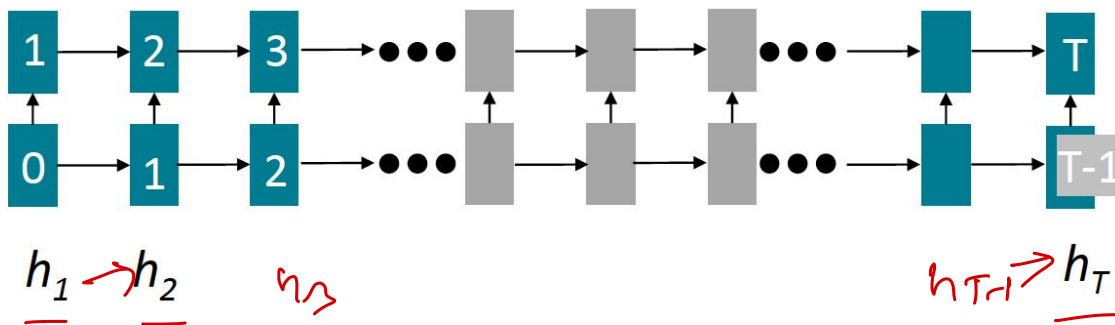


*The* **chef** *who ...*  **was**

Info of ~~chef~~ has gone through O(sequence length) many layers!

https://web.stanford.edu/class/cs224n/

# Issues with recurrent models

## *Lack of parellilizability*

- Future RNN hidden states can't be computed in full before past RNN hidden states have been computed
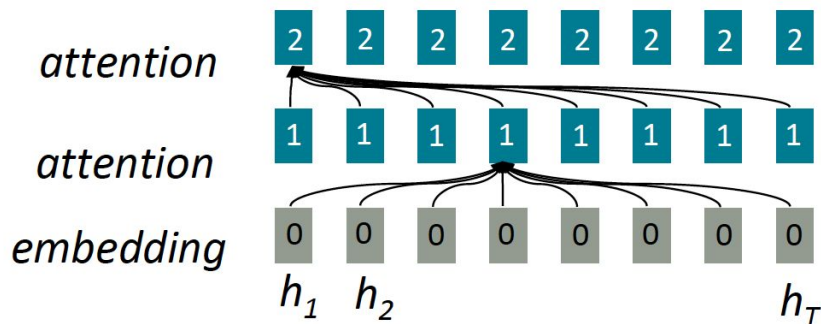


Inhibits training on very large datasets!

Numbers indicate min # of steps before a state can be computed

https://web.stanford.edu/class/cs224n/

# If not recurrence, then what?

## Attention

- Given a word as query, attention can be used to access and incorporate information from a set of values (other words)
- Can we do this within a single sentence?
- All words can interact with each other and computation can be done in parallel!!



*attention*

*attention*

*embedding*

$h_1$  $h_2$     $h_T$

All words attend to all words in previous layer; most arrows here are omitted

https://web.stanford.edu/class/cs224n/

# Transformer Encoder-Decoder

- Transformer is introduced as an **encoder-decoder architecture**; later we will see **encoder-only** & **decoder-only** transformers
- **Encoder** produces a sophisticated representation of the source sequence that the decoder will use to condition its generation process
- **Decoder** generates one token at the time to produce a target sequence; it produces representations that combine the history and a new token
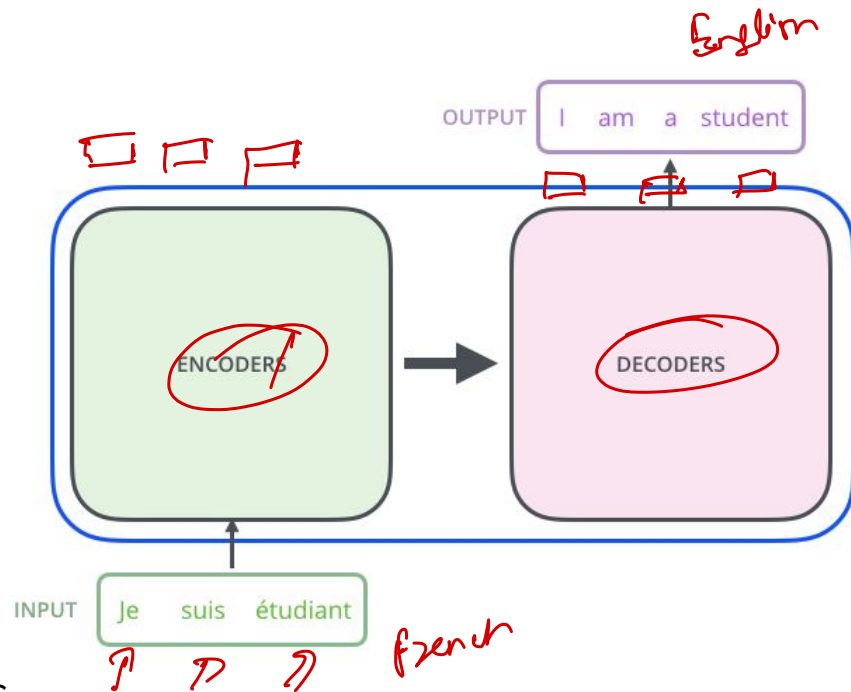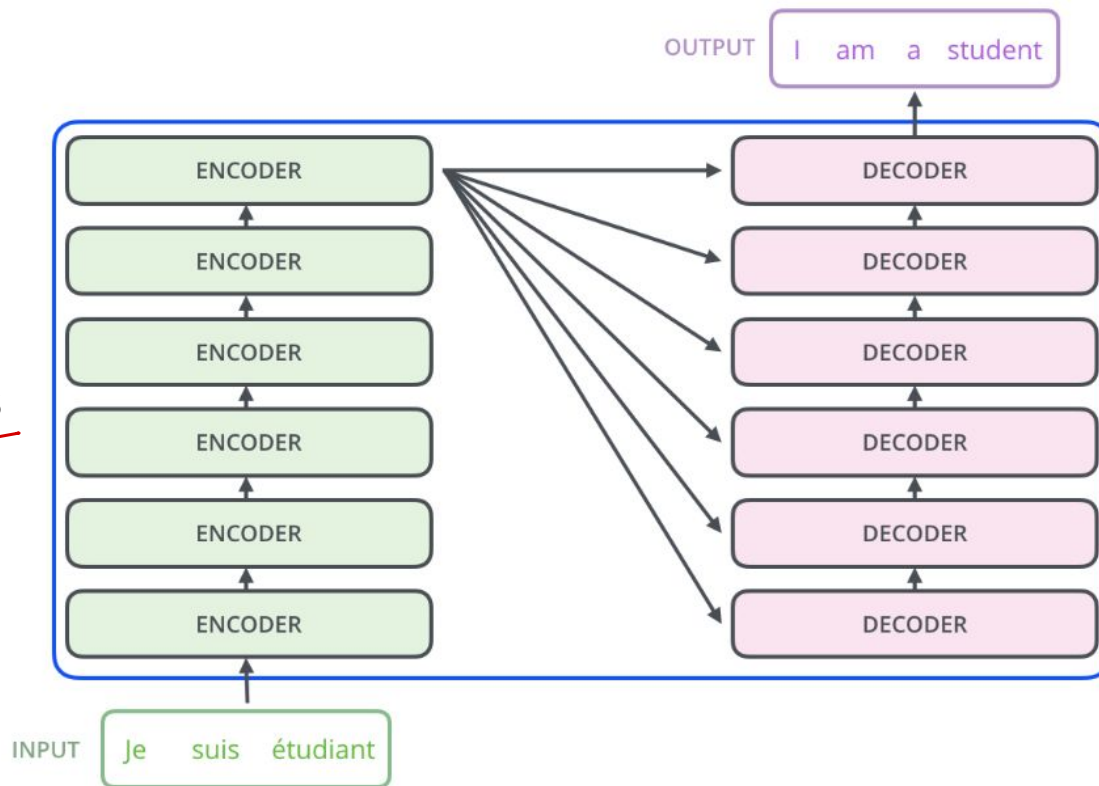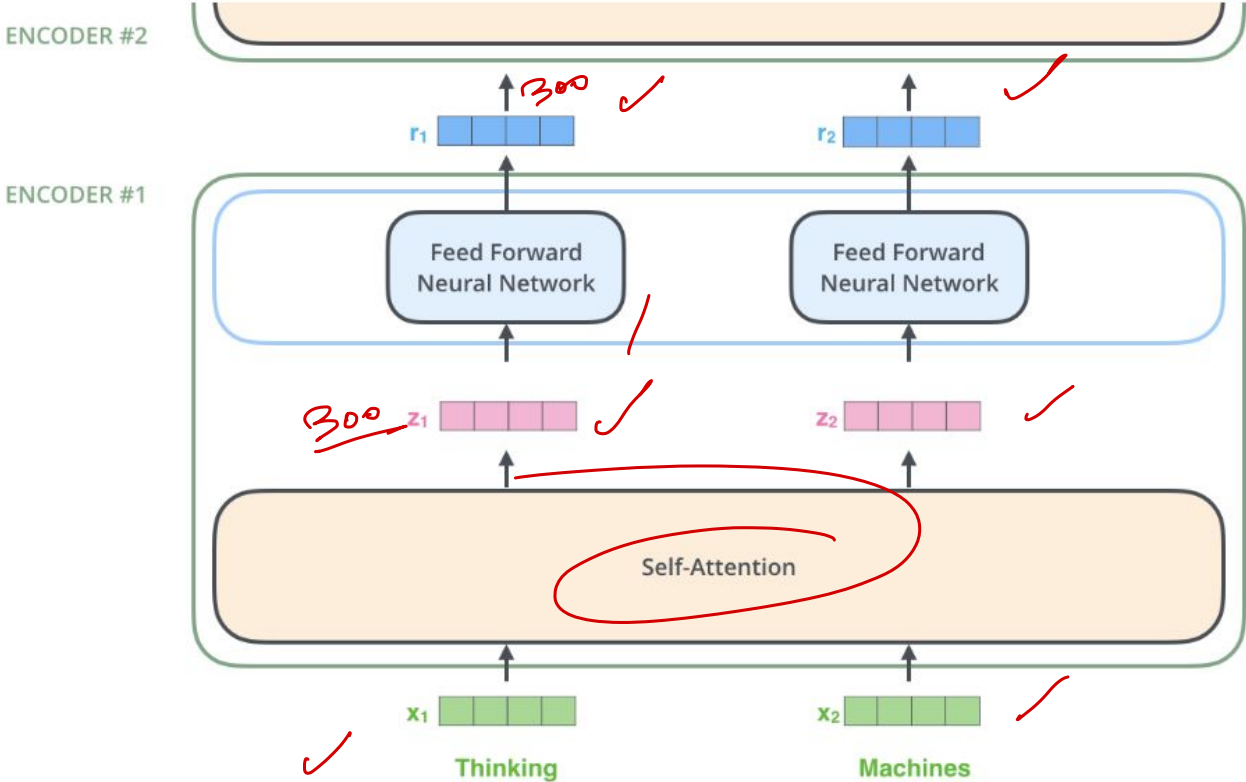
Figure: Jay Alammar

OUTPUT: I am a student

ENCODER
ENCODER
ENCODER
ENCODER
ENCODER
ENCODER

DECODER
DECODER
DECODER
DECODER
DECODER
DECODER

A stack of encoder blocks

A stack of decoder blocks

INPUT: Je suis étudiant

Figure: Jay Alammar

# Each encoder block consists of **self-attention** & **FFNN**



Figure: Jay Alammar

Deeper layers get outputs of the previous layers as inputs

Input to each encoder block has the same size as original token embeddings

In the first layer, inputs are static token emeddings

# Intuition for Self-Attention

## *Problem with static embeddings (word2vec)*

They are static!  The embedding for a word doesn't reflect how its meaning changes in context.

The chicken didn't cross the road because it was too tired

What is the meaning represented in

the static embedding for "it"?

https://web.stanford.edu/~jurafsky/slp3/

# Intuition for attention

The chicken didn't cross the road because it

## What should be the properties of "it"?

The chicken didn't cross the road because it was too **tired**
The chicken didn't cross the road because it was too **wide**

At this point in the sentence, it's probably referring to either the animal or the street

*Contextual representation*
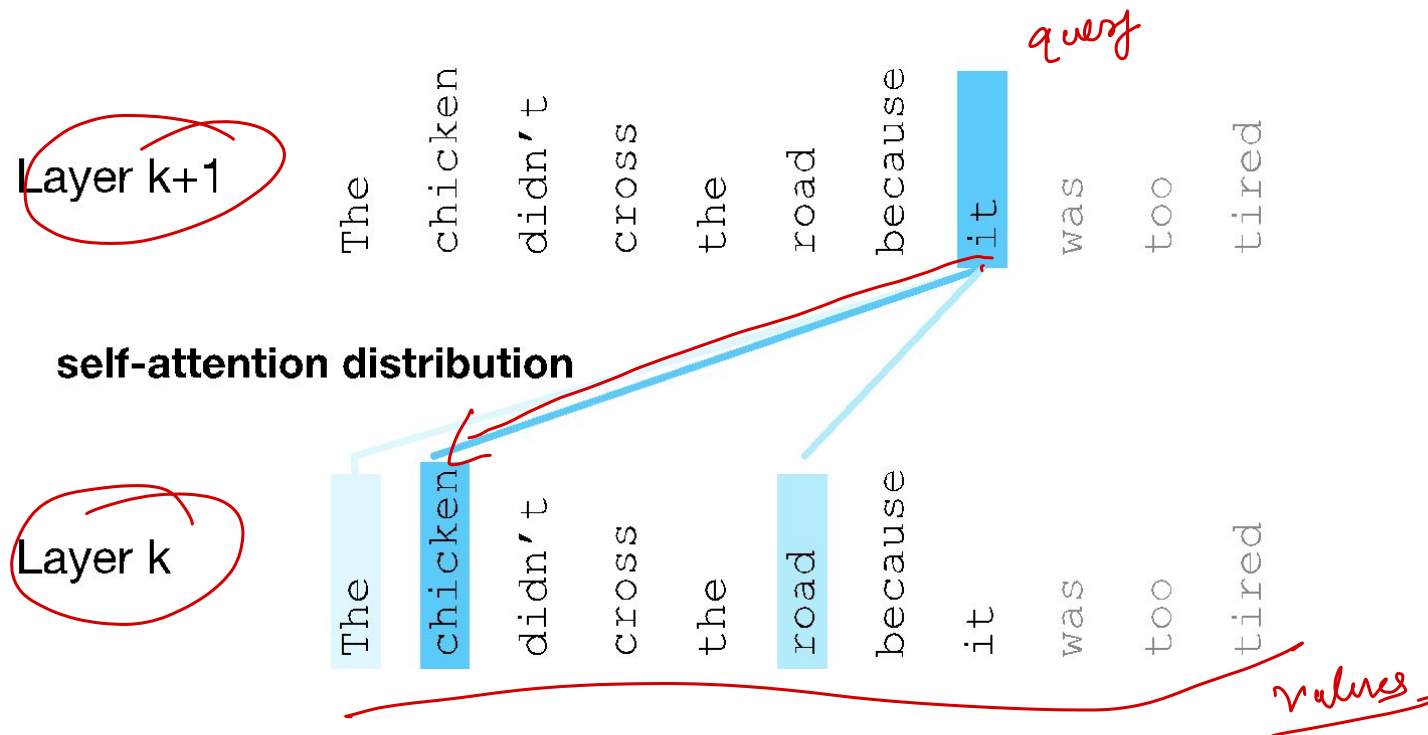
https://web.stanford.edu/~jurafsky/slp3/

# Intuition of attention

Build up the contextual embedding from a word by selectively integrating information from all the neighboring words

We say that a word "attends to" some neighboring words more than others

# Intuition of attention
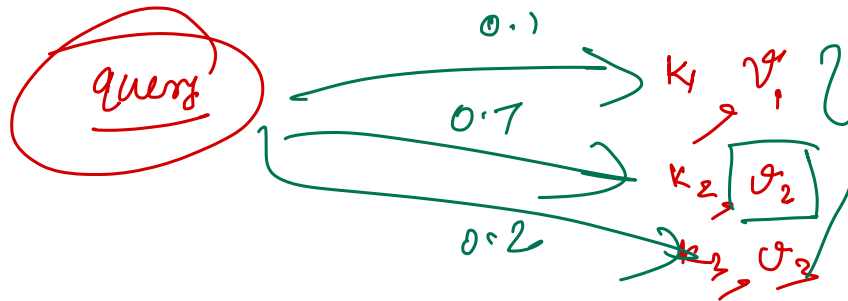
columns corresponding to input tokens



Layer k+1

query

self-attention distribution

Layer k

values

https://web.stanford.edu/~jurafsky/slp3/
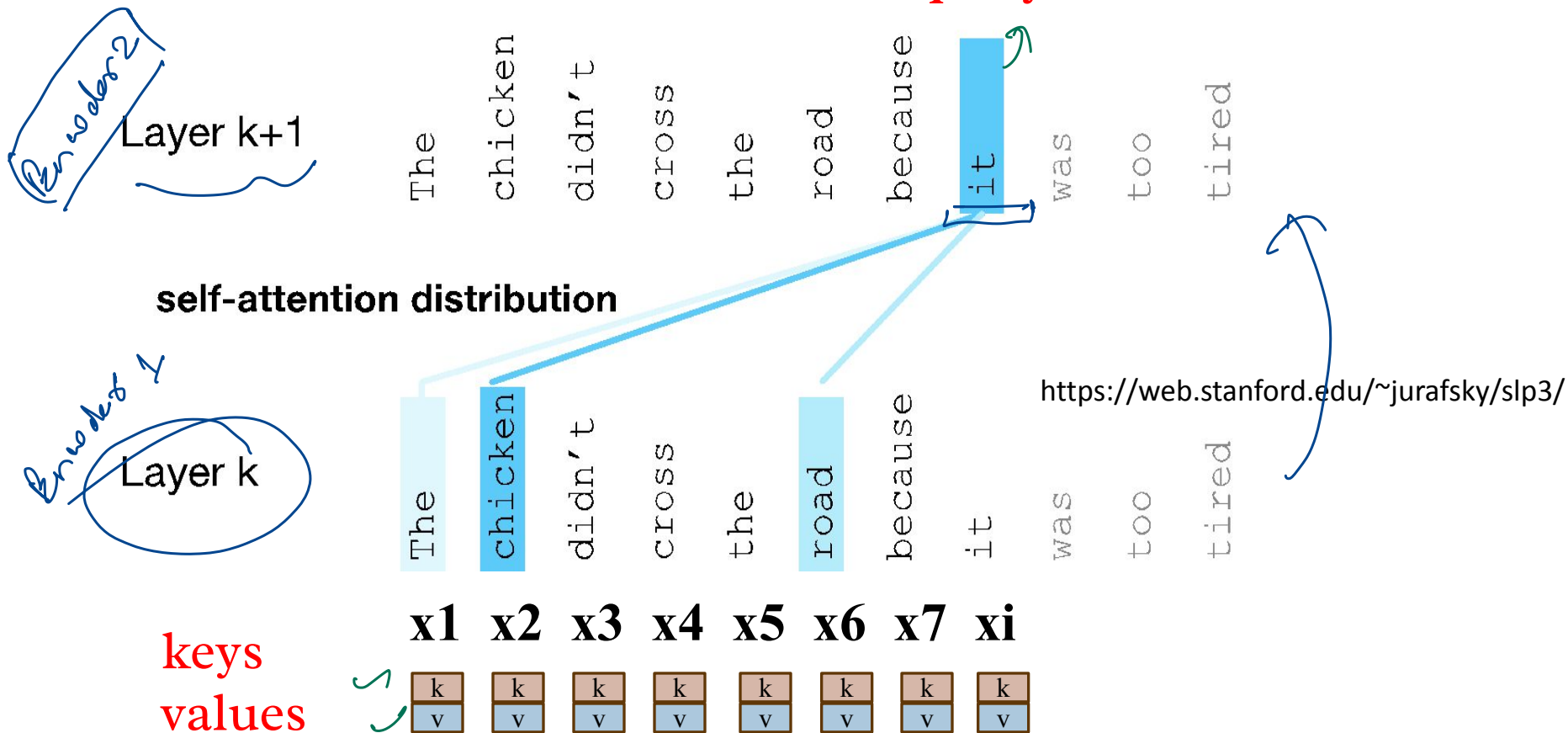
# Intuition for Self-attention

*Attention is based on key/value/query concept -- analogous to retrieval systems.*

**When you search for videos on Youtube**

- The search engine will map your query (text in the search bar) against a set of keys (video title, description, etc.) associated with candidate videos in their database

- It will then present you the best matched videos (values).

# Intuition of attention:



**query**

Layer k+1

The chicken didn't cross the road because it was too tired

**self-attention distribution**

https://web.stanford.edu/~jurafsky/slp3/

Layer k

The chicken didn't cross the road because it was too tired

**x1  x2  x3  x4  x5  x6  x7  xi**

**keys**
**values**

| k | k | k | k | k | k | k | k |
|---|---|---|---|---|---|---|---|
| v | v | v | v | v | v | v | v |

Encoder 2

Encoder 1

# An Actual Attention Head: slightly more complicated

We'll use matrices to project each vector $\mathbf{x}_i$ into a representation of its role as query, key, value:

- **query**: $W^Q$
- **key**: $W^K$
- **value**: $W^V$

3 learnable parameter (matrices)

$$\mathbf{q}_i = \mathbf{x}_i W^Q; \quad \mathbf{k}_i = \mathbf{x}_i W^K; \quad \mathbf{v}_i = \mathbf{x}_i W^V$$

https://web.stanford.edu/~jurafsky/slp3/

# An Actual Attention Head: slightly more complicated

Given these 3 representation of $x_i$

To compute similarity of current element $x_i$ with some prior element $x_j$

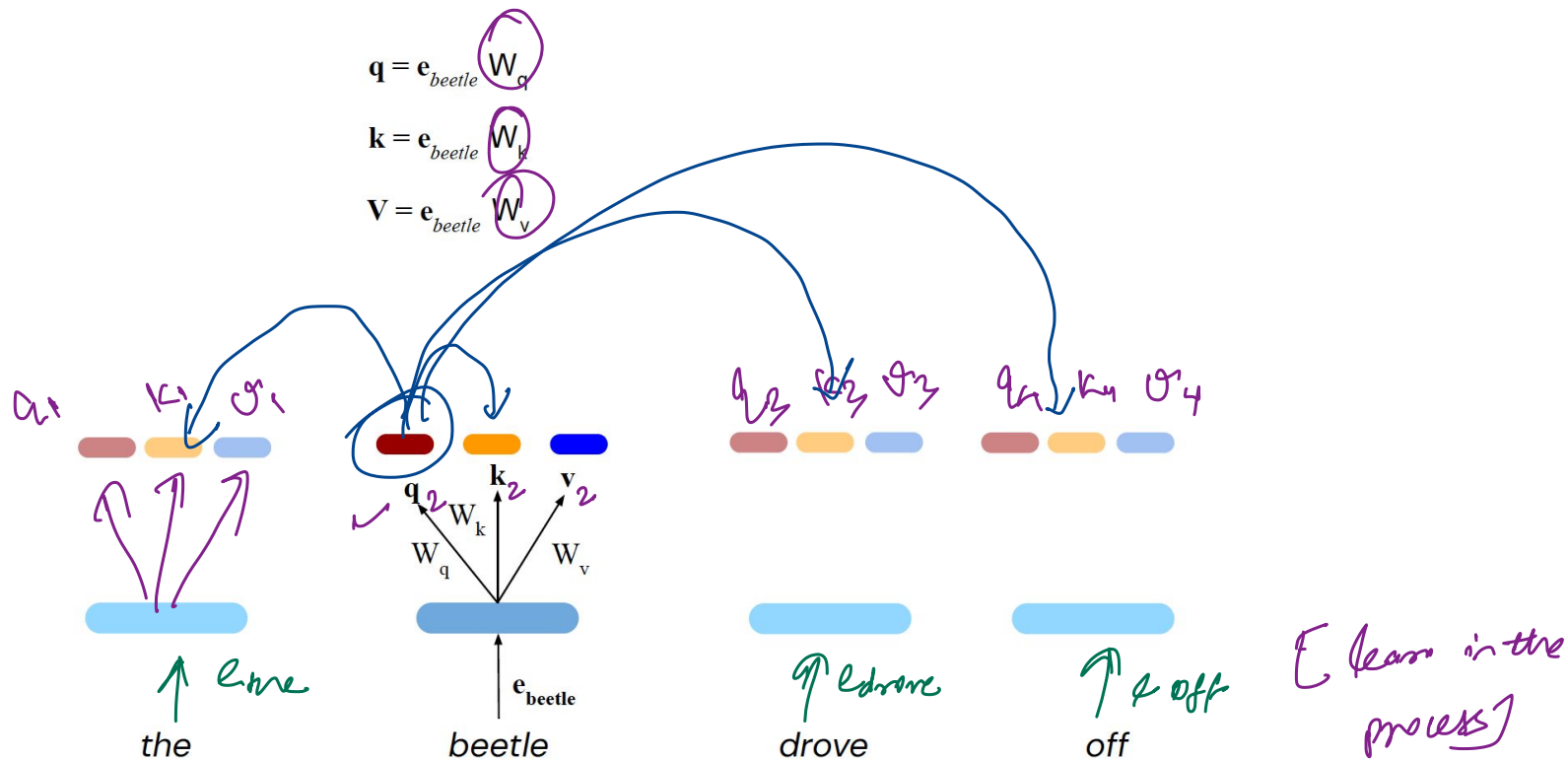We'll use dot product between $q_i$ and $k_j$.

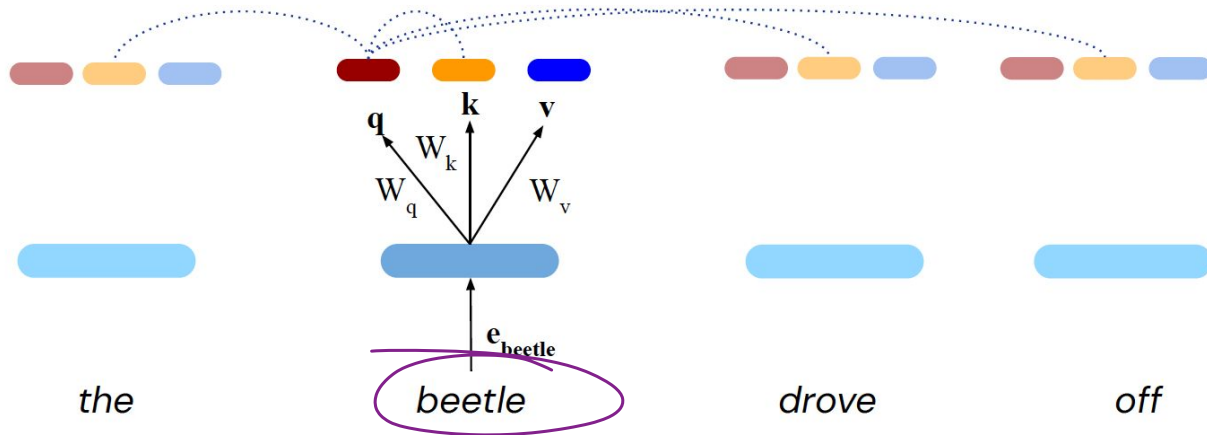And instead of summing up $x_j$, we'll sum up $v_j$

$$\mathbf{q}_i = \mathbf{x}_i \mathbf{W}^Q; \quad \mathbf{k}_i = \mathbf{x}_i \mathbf{W}^K; \quad \mathbf{v}_i = \mathbf{x}_i \mathbf{W}^V$$

# Transformers: Self-attention over input

$$\mathbf{q} = \mathbf{e}_{beetle} W_q$$

$$\mathbf{k} = \mathbf{e}_{beetle} W_k$$

$$\mathbf{V} = \mathbf{e}_{beetle} W_v$$



$q_1$   $k_1$   $v_1$

$t$

$q_3$   $k_3$   $v_3$      $q_4$   $k_4$   $v_4$

$\mathbf{q}_2$   $\mathbf{k}_2$   $\mathbf{v}_2$

$W_q$   $W_k$   $W_v$

$\mathbf{e}_{beetle}$

*the*       *beetle*       *drove*       *off*

ⅇᵢⅈₑ       q drove       q off

[ learn in the process ]

# Self-attention over input embeddings



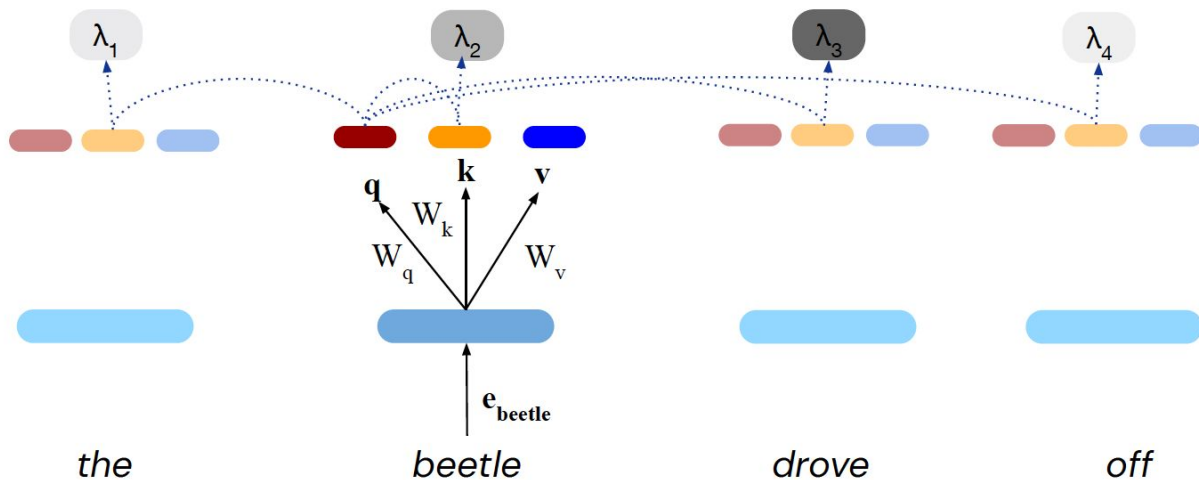Source: *https://deepmind.com/learning-resources/deep-learning-lecture-series-2020*
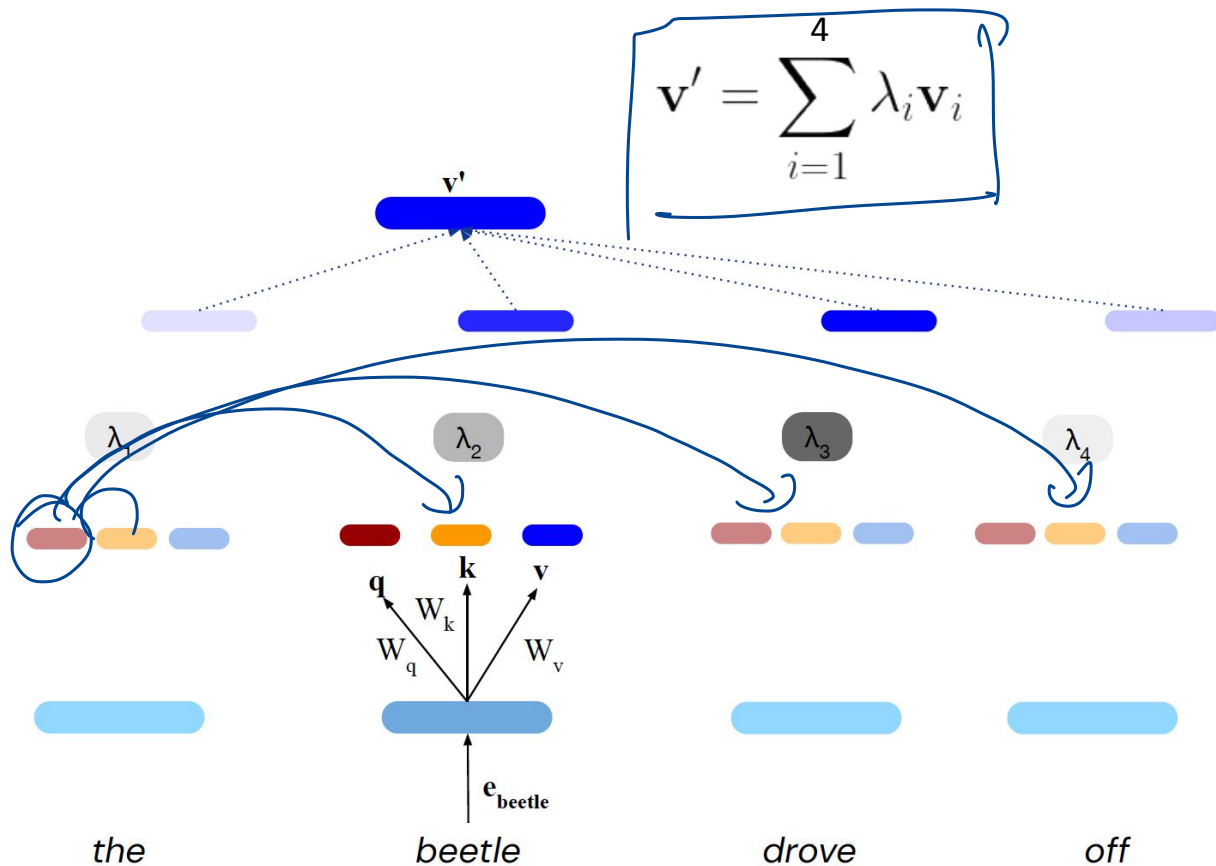
# Self-attention over input embeddings

$$\lambda_i = \frac{e^{\mathbf{q} \cdot \mathbf{k_i}}}{\sum_{i=1}^{4} e^{\mathbf{q} \cdot \mathbf{k_i}}}$$
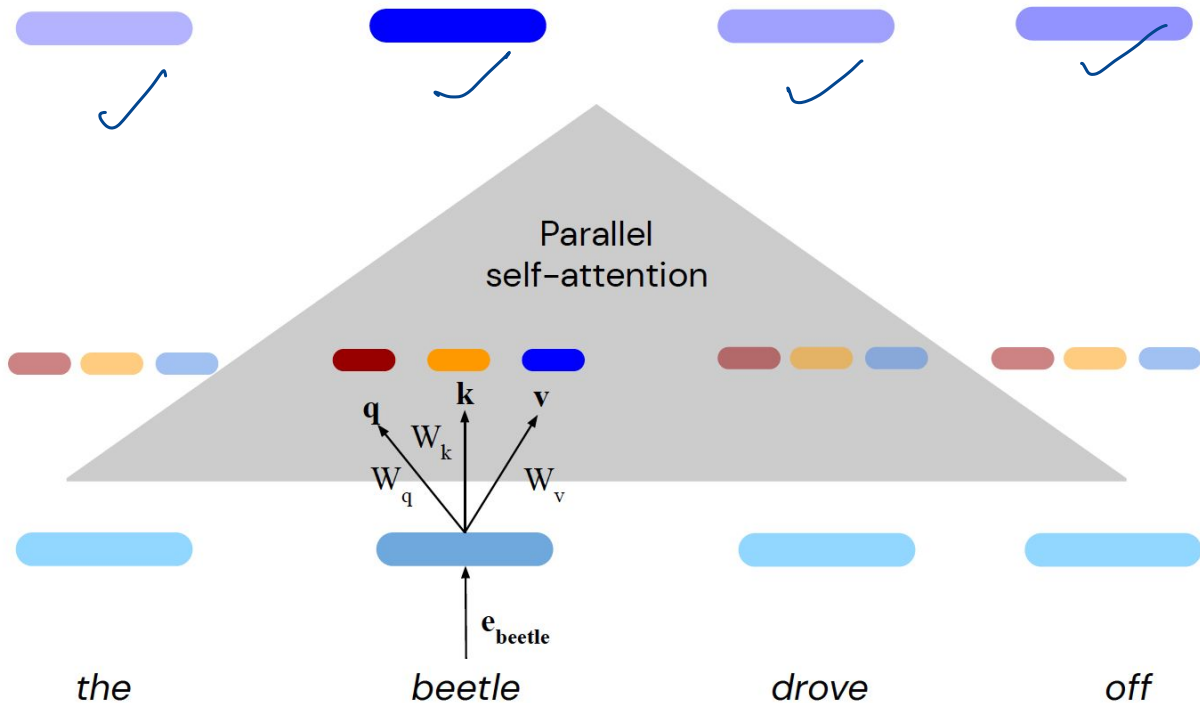
$$\Sigma \, \lambda_i \, \vartheta_i$$

# Self-attention over input embeddings



$$\mathbf{v}' = \sum_{i=1}^{4} \lambda_i \mathbf{v}_i$$

Source: *https://deepmind.com/learning-resources/deep-learning-lecture-series-2020*

# Self-attention over all words (in parallel)



Parallel self-attention

$\mathbf{q}$  $\mathbf{k}$  $\mathbf{v}$

$W_k$

$W_q$  $W_v$

$\mathbf{e_{beetle}}$

*the*  *beetle*  *drove*  *off*

Source: *https://deepmind.com/learning-resources/deep-learning-lecture-series-2020*

# Self-attention over all words (in parallel)

# Self-attention: In equations

hardmaru ✔
@hardmaru · Follow

The most important formula in deep learning after 2018

**Self-Attention**

**What is self-attention?** Self-attention calculates a weighted average of feature representations with the weight proportional to a similarity score between pairs of representations. Formally, an input sequence of $n$ tokens of dimensions $d$, $X \in \mathbf{R}^{n \times d}$, is projected using three matrices $W_Q \in \mathbf{R}^{d \times d_q}$, $W_K \in \mathbf{R}^{d \times d_k}$, and $W_V \in \mathbf{R}^{d \times d_v}$ to extract feature representations $Q$, $K$, and $V$, referred to as query, key, and value respectively with $d_k = d_q$. The outputs $Q$, $K$, $V$ are computed as

$$Q = XW_Q, \quad K = XW_K, \quad V = XW_V. \quad (1)$$

So, self-attention can be written as,

$$S = D(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_q}}\right)V, \quad (2)$$

where softmax denotes a *row-wise* softmax normalization function. Thus, each element in $S$ depends on all other elements in the same row.
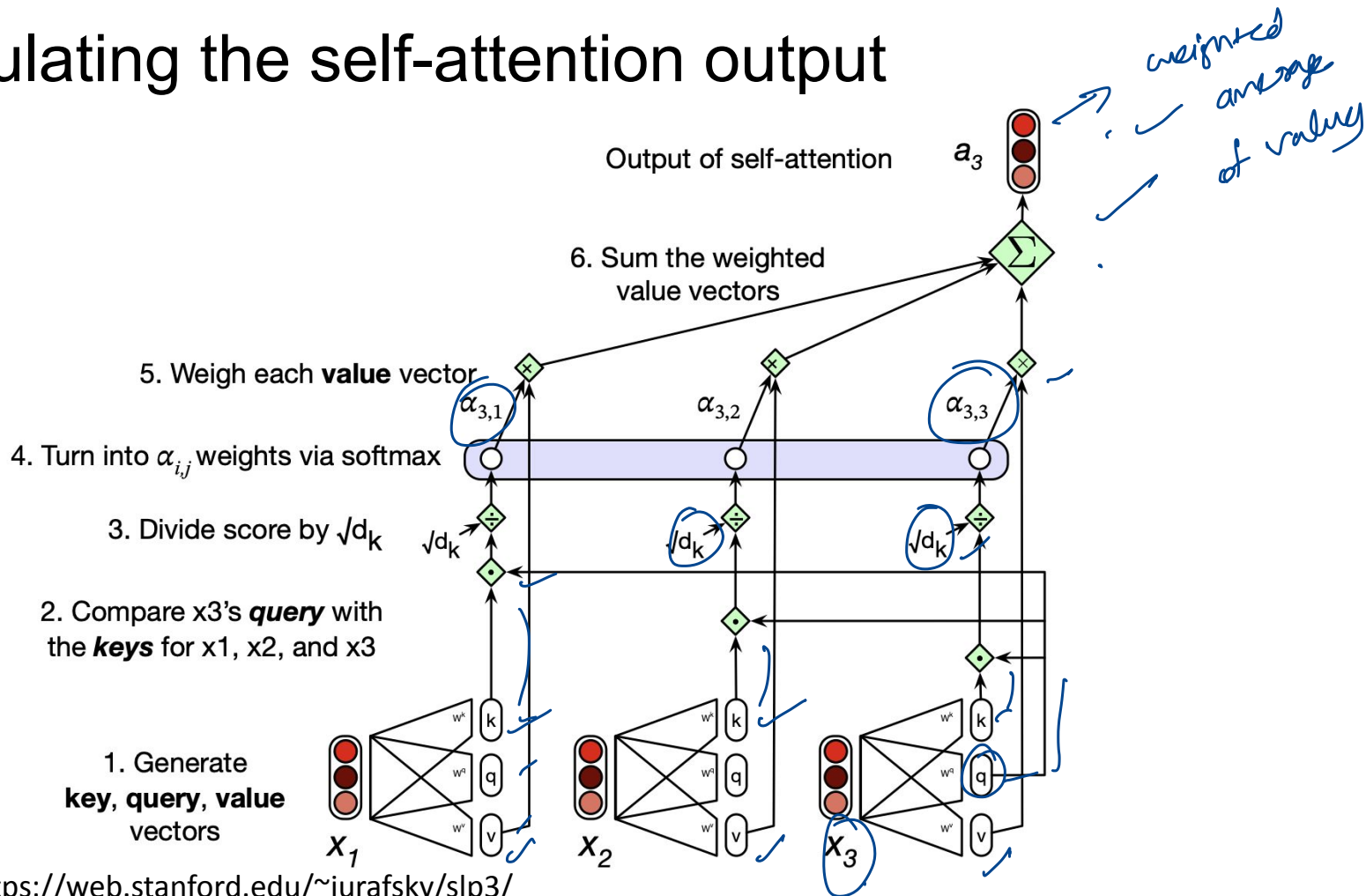
7:38 AM · Feb 10, 2021

❤ 3.1K    💬 Reply    🔗 Copy link

$$\mathbf{q}_i = \mathbf{x}_i \mathbf{W}^Q; \quad \mathbf{k}_j = \mathbf{x}_j \mathbf{W}^K; \quad \mathbf{v}_j = \mathbf{x}_j \mathbf{W}^V$$

$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{d_k}}$$

$$\alpha_{ij} = \text{softmax}(\text{score}(\mathbf{x}_i, \mathbf{x}_j))$$

$$\mathbf{a}_i = \sum \alpha_{ij} \mathbf{v}_j$$

**Scaled dot-product**: *more on this later*

Source: *https://theaisummer.com/self-attention/* *https://web.stanford.edu/~jurafsky/slp3/*

# Calculating the self-attention output



Output of self-attention   $a_3$

weighted average of values

6. Sum the weighted value vectors

5. Weigh each **value** vector   $\alpha_{3,1}$   $\alpha_{3,2}$   $\alpha_{3,3}$

4. Turn into $\alpha_{i,j}$ weights via softmax

3. Divide score by $\sqrt{d_k}$   $\sqrt{d_k}$   $\sqrt{d_k}$   $\sqrt{d_k}$

2. Compare x3's **query** with the **keys** for x1, x2, and x3

1. Generate **key**, **query**, **value** vectors   $x_1$   $x_2$   $x_3$

https://web.stanford.edu/~jurafsky/slp3/

# Try this problem

Suppose, you give the following input to your transformer encoder: {flying, arrows} The input embeddings for these two words are **[0,1,1,1,1,0] and [1,1,0,-1,-1,1]**, respectively. Suppose you are trying to represent the first word 'flying' with the help of self-attention in the first encoder. For the first attention head, the query, key and value matrices just take the 2 dimensions from the input each. Thus, the first 2 dimensions define the query vector, and so on. What will be the self-attention output for the word 'flying' corresponding to this attention head. You are using the scaled dot vector.