# Transformer
## Attention Is All You Need

Seyedi

# Attention Is All You Need

**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

[†]Work performed while at Google Brain.
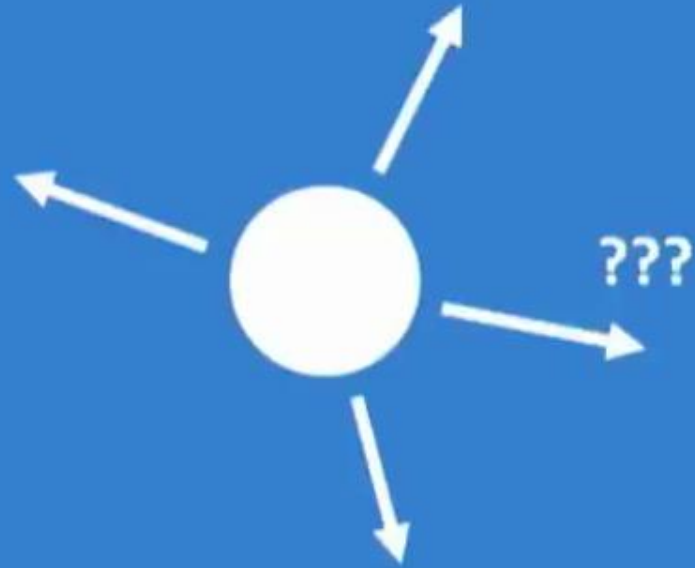[‡]Work performed while at Google Research.

# Sequence

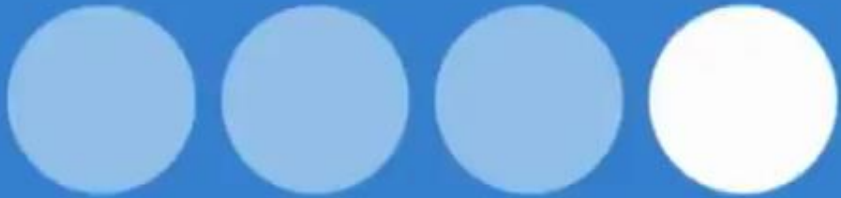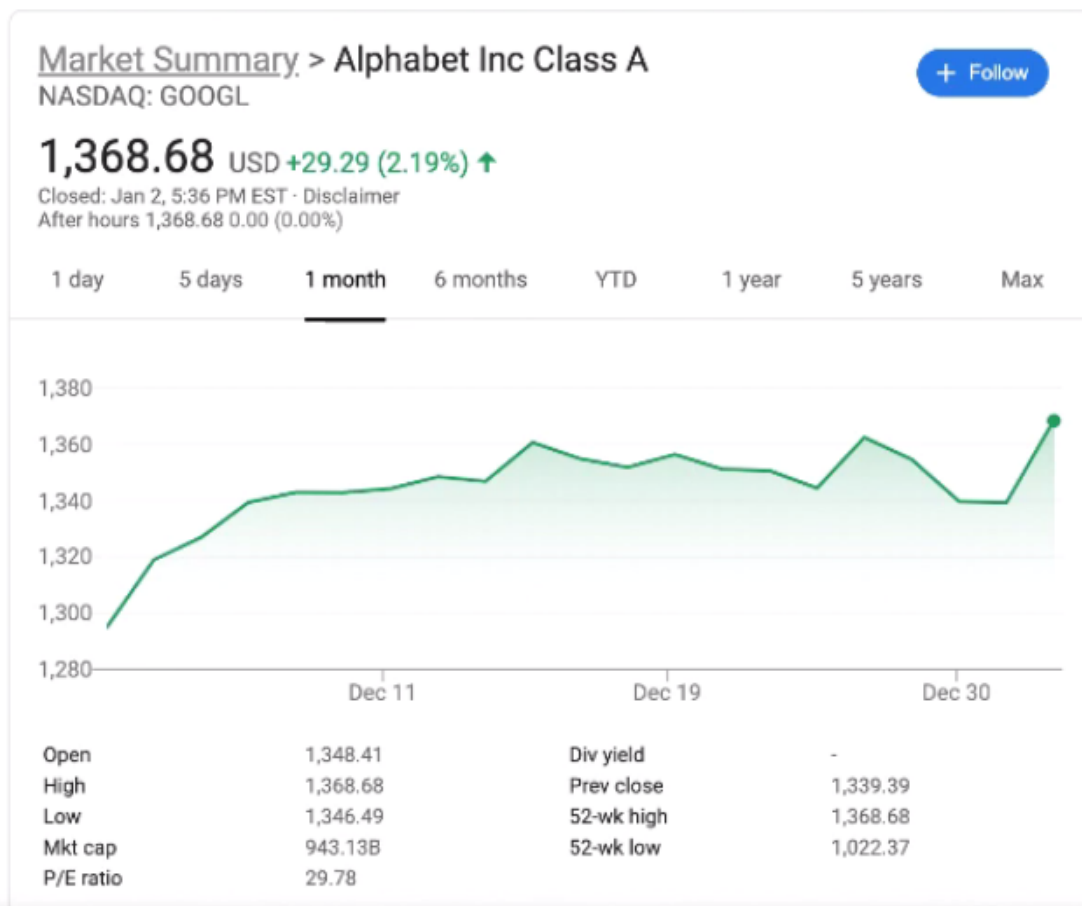Given an image of a ball,
can you predict where it will go next?

???

# Sequence

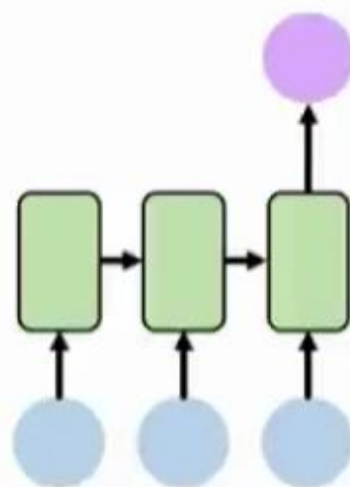The unicorn is scotland's national animal

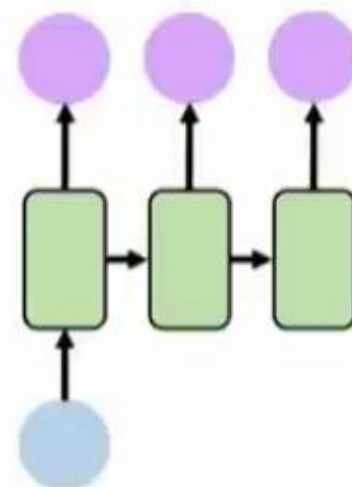Sequences

# Sequence



## Sequence Modeling Applications

**One to One**
**Binary Classification**
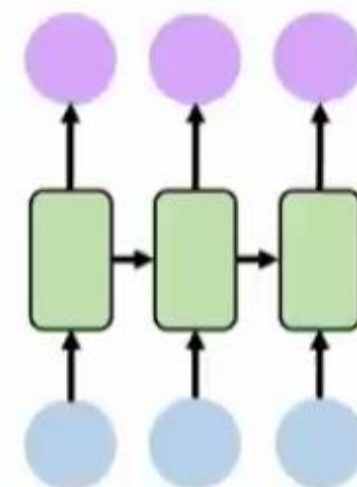
"Will I pass this class?"
Student → Pass?

**Many to One**
**Sentiment Classification**

Ivar Hagendoorn

The @MIT Introduction to #DeepLearning is definitely one of the best courses of its kind currently available online
introtodeeplearning.com

**One to Many**
**Image Captioning**

"A baseball player throws a ball."

**Many to Many**
**Machine Translation**

# Handling Individual Time Steps



$$\hat{y}_t = f(x_t)$$

# Neurons with Recurrence



$$\hat{y}_t = f(x_t, h_{t-1})$$

output      input    past memory

# Recurrent Neural Networks

# Recurrent Neural Networks

Recurrent Neural Networks has a short reference window

As aliens entered our planet          and began to colonize earth a certain group of extraterrestrials ...
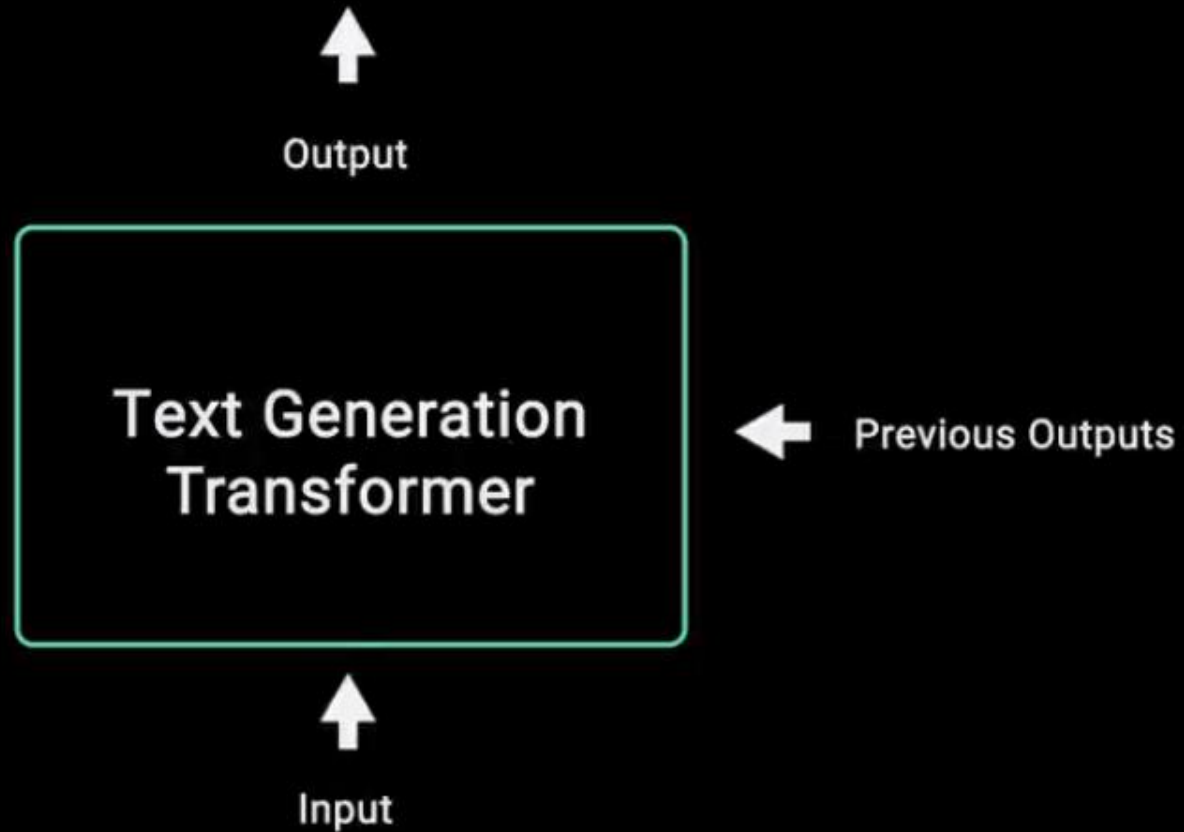
# GRUs and LSTMs

GRU's and LSTM's have a longer reference window than RNN's

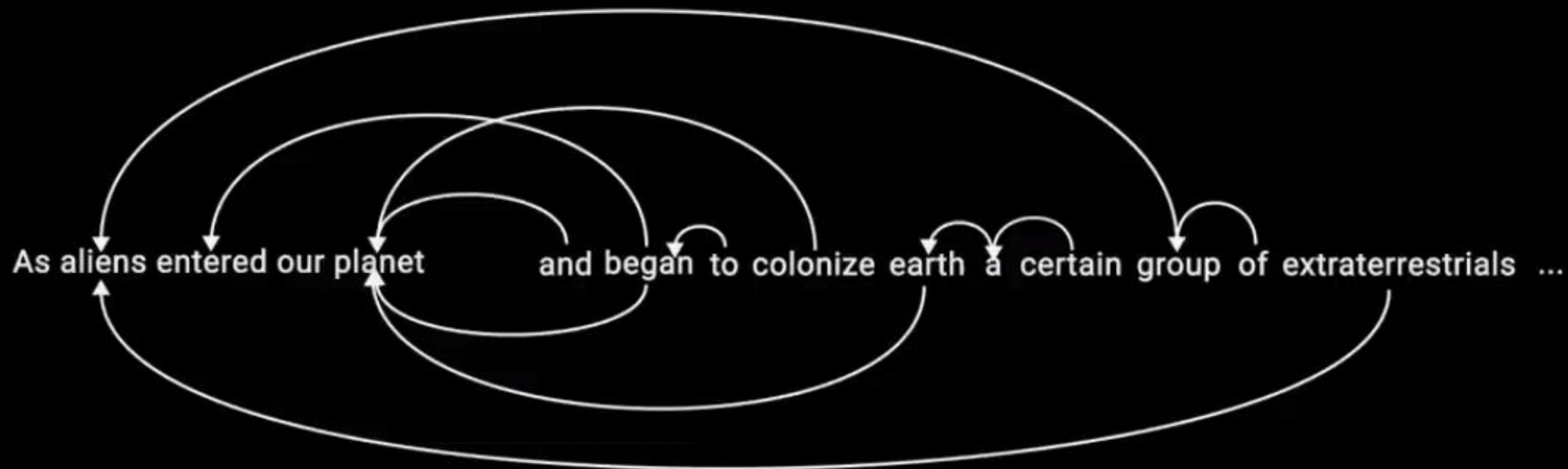As aliens entered our planet | and began to colonize earth a certain group of extraterrestrials ...

# Text Generation and Attention
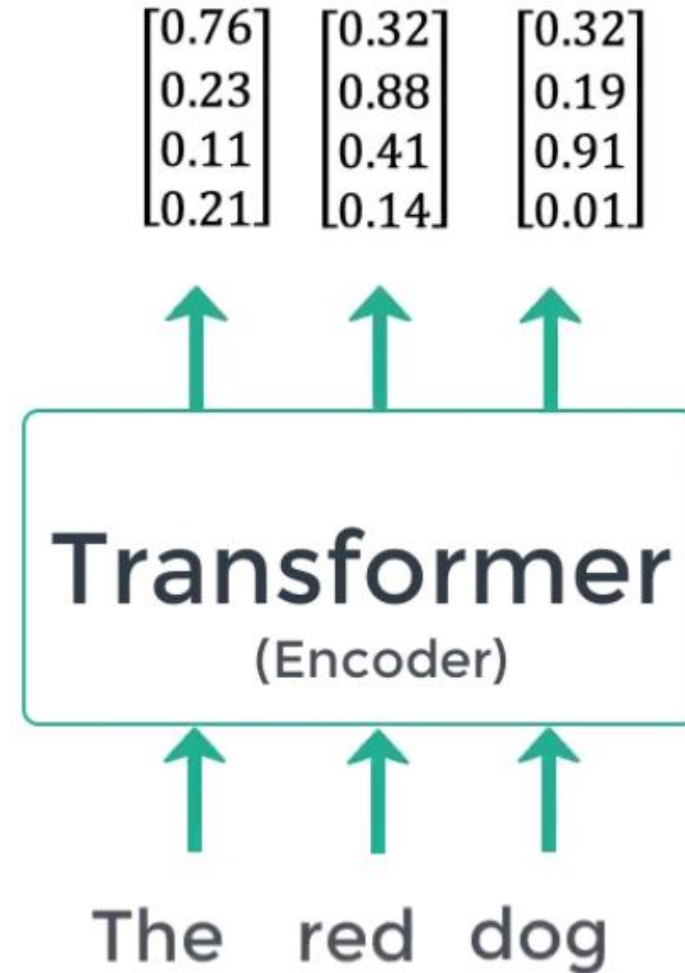
# Text Generation and Attention

# Text Generation and Attention



As aliens entered our planet     and began to colonize earth a certain group of extraterrestrials ...
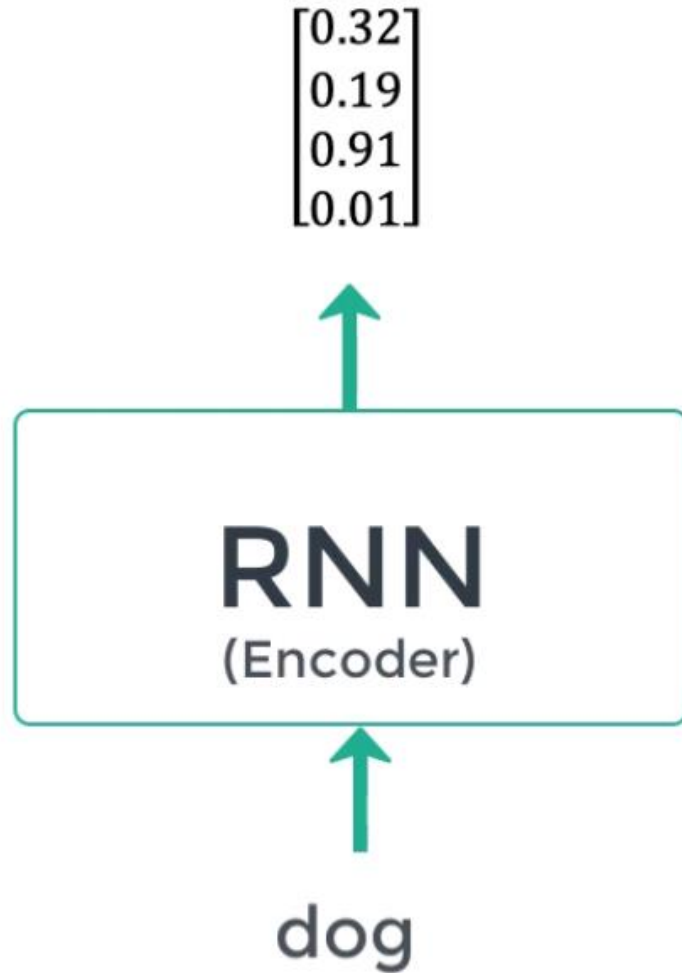
# Attention Mechanism has an infitnite reference window

As aliens entered our planet       and began to colonize earth a certain group of extraterrestrials ...

# RNNs vs Transformers

$$\begin{bmatrix} 0.32 \\ 0.19 \\ 0.91 \\ 0.01 \end{bmatrix}$$

$$\begin{bmatrix} 0.76 \\ 0.23 \\ 0.11 \\ 0.21 \end{bmatrix} \begin{bmatrix} 0.32 \\ 0.88 \\ 0.41 \\ 0.14 \end{bmatrix} \begin{bmatrix} 0.32 \\ 0.19 \\ 0.91 \\ 0.01 \end{bmatrix}$$

**RNN**
(Encoder)

**Transformer**
(Encoder)

dog

The   red   dog

# RNNs vs Transformers

## Challenges with RNNs

- Long range dependencies

- Gradient vanishing and explosion

- Large # of training steps

- Recurrence prevents parallel computation

## Transformer Networks

- Facilitate long range dependencies

- No gradient vanishing and explosion

- Fewer training steps

- No recurrence that facilitate parallel computation

Transformers
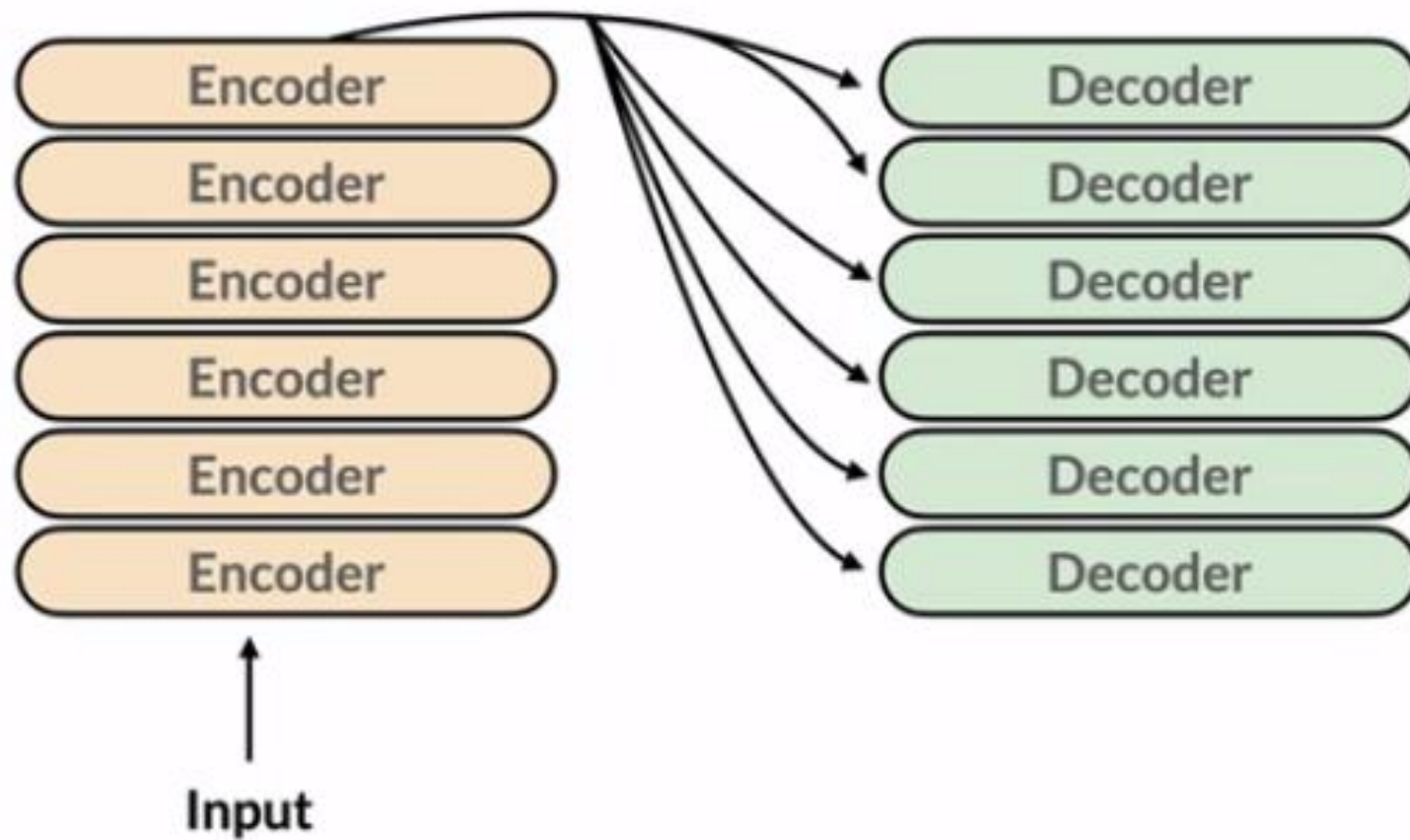Encoder

Transformers
Decoder

# Transformer Network

# Transformer Network

# Transformer Network

Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Add & Norm

Masked
Multi-Head
Attention

N×

N×

Positional
Encoding

Positional
Encoding

Input
Embedding

Output
Embedding

Inputs

Outputs
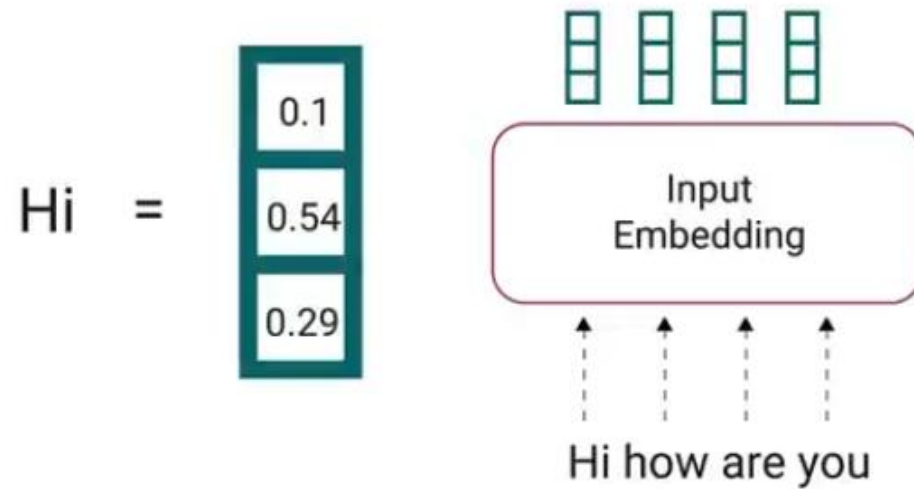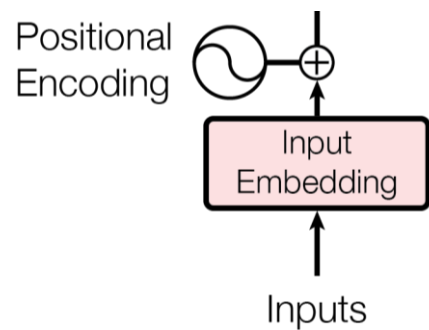(shifted right)

# Encoder

# Embedding & Positional Encoding

# Embedding & Positional Encoding

## 1. Input Embedding

Hi = 
| 0.1 |
| 0.54 |
| 0.29 |

Input Embedding

Hi how are you

Positional Encoding ⊕

Input Embedding

Inputs

# Embedding & Positional Encoding

| Positional Encoding | ▯ | ▯ | ▯ | ▯ |
|---|---|---|---|---|
| Time Step | 1 | 2 | 3 | 4 |

Positional Encoding ⟳—⊕
Input Embedding
Inputs

# Embedding & Positional Encoding

Positional Encoding

Time Step     1        2        3        4

$$PE(pos, 2i + 1) = cos(\frac{pos}{10000^{2i/dmodel}})$$

$$PE(pos, 2i) = sin(\frac{pos}{10000^{2i/dmodel}})$$
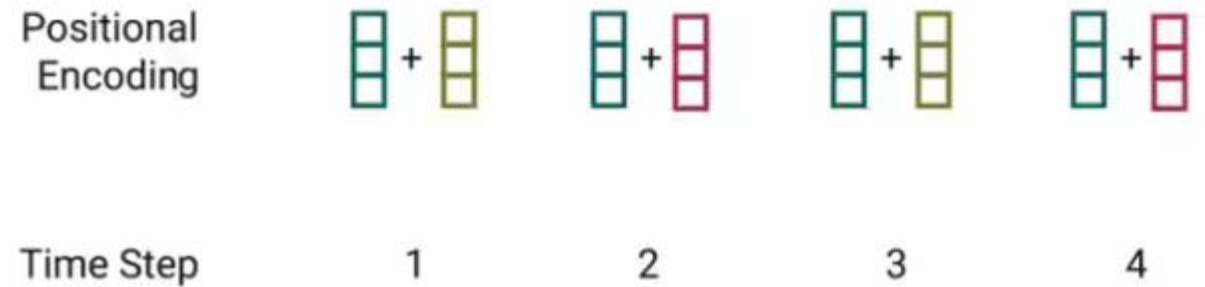
Positional Encoding

Input Embedding

Inputs

# Embedding & Positional Encoding

## 2. Positional Encoding

Positional Input Embeddings

Positional Encoding

Time Step      1          2          3          4

$$PE(pos, 2i+1) = cos(\frac{pos}{10000^{2i/dmodel}})$$

$$PE(pos, 2i) = sin(\frac{pos}{10000^{2i/dmodel}})$$

Positional Encoding

Input Embedding
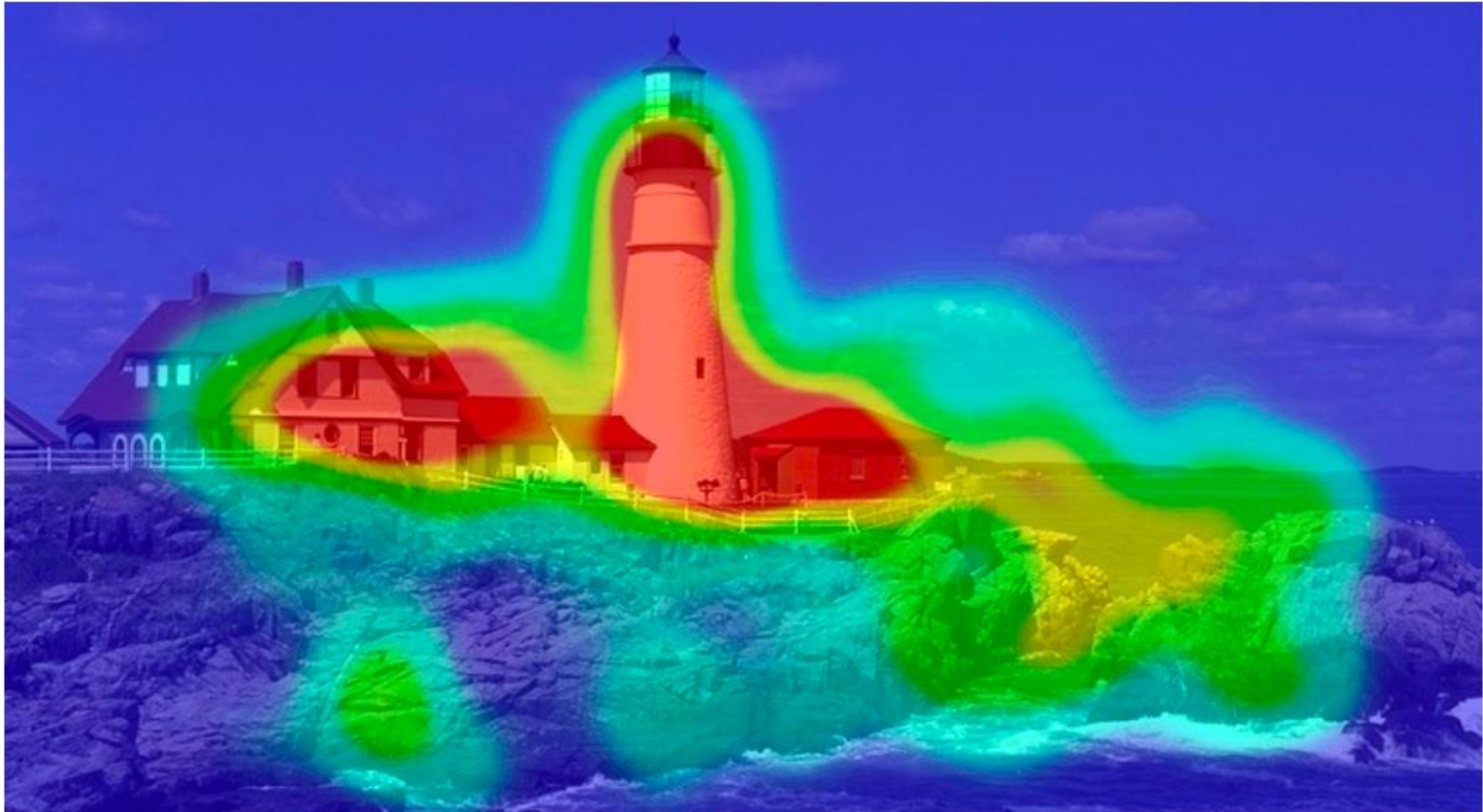
Inputs

# Multi-Head Self-Attention

# Attention

- Attention in Computer Vision
  - 2014: Attention used to highlight important parts of an image that contribute to a desired output



- Attention in NLP
  - 2015: Aligned machine translation
  - 2017: Language modeling with **Transformer networks**

# Attention

# Attention

# Attention

# Intuition behind Self-Attention



Attending to the most important parts of an input.

# Intuition behind Self-Attention



Attending to the most important parts of an input.

1. Identify which parts to attend to
2. Extract the features with high attention

# Intuition behind Self-Attention



Attending to the most important parts of an input.

1. Identify which parts to attend to ← Similar to a search problem!
2. Extract the features with high attention

# Understanding Attention with Search

# Understanding Attention with Search

# Understanding Attention with Search



**Query (Q)**

**Key (K₁)** — How similar is the key to the query?

**Key (K₂)**

**Key (K₃)**

I. **Compute attention mask:** how similar is each key to the desired query?

# Understanding Attention with Search



**Query (Q)**

**Key (K₂)**

**Key (K₃)**

How similar is the key to the query?
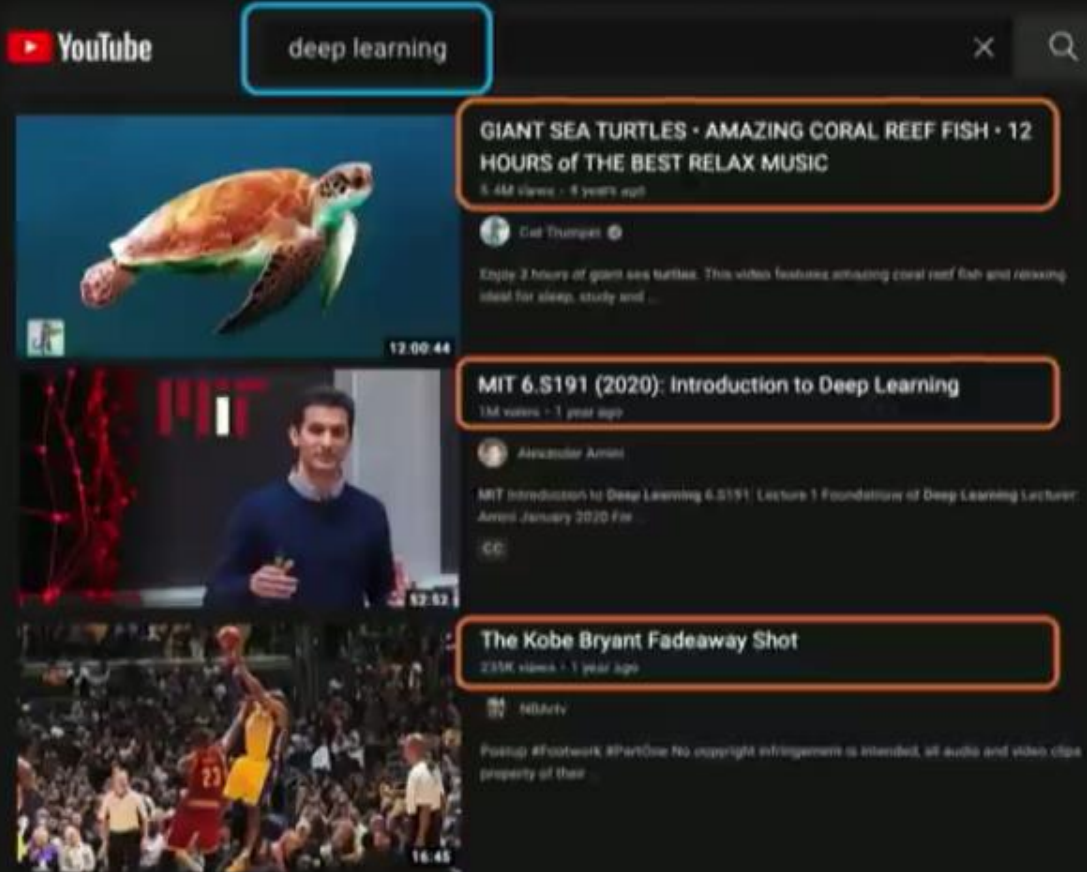
1. **Compute attention mask:** how similar is each key to the desired query?

# Understanding Attention with Search



Query (Q)

Key (K₂)

Value (V)

2. **Extract values based on attention:**
Return the values highest attention

# Learning Self-Attention

Goal: identify and attend to most important features in input.

1. Encode **position** information

2. Extract query, key, value for search

3. Compute attention weighting

4. Extract features with high attention



Data is fed in all at once! Need to encode position information to understand order.

# Learning Self-Attention





Multi-Head
Attention

# Learning Self-Attention



query     key     value

Linear     Linear     Linear

Multi-Head Attention

# Learning Self-Attention

Goal: identify and attend to most important features in input.

1. Encode **position** information
2. Extract **query, key, value** for search
3. Compute attention weighting
4. Extract features with high attention

# Learning Self-Attention

**Goal:** identify and attend to most important features in input.

1. Encode **position** information
2. Extract **query, key, value** for search
3. Compute **attention weighting**
4. Extract features with high attention

**Attention score:** compute pairwise similarity between each query and key

How to compute similarity between two sets of features?



Also known as the "cosine similarity"
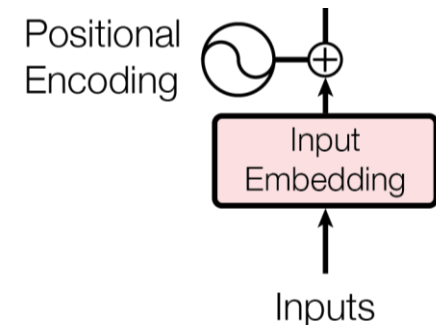
# Learning Self-Attention

# Learning Self-Attention

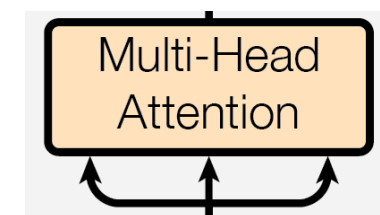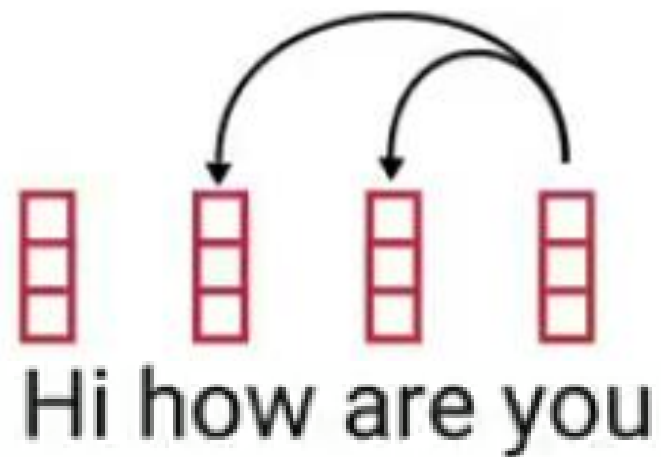**Goal:** identify and attend to most important features in input.

1. Encode **position** information
2. Extract **query, key, value** for search
3. Compute **attention weighting**
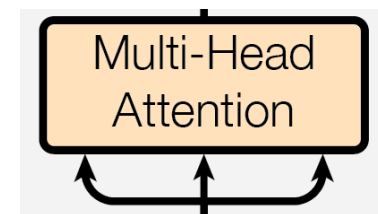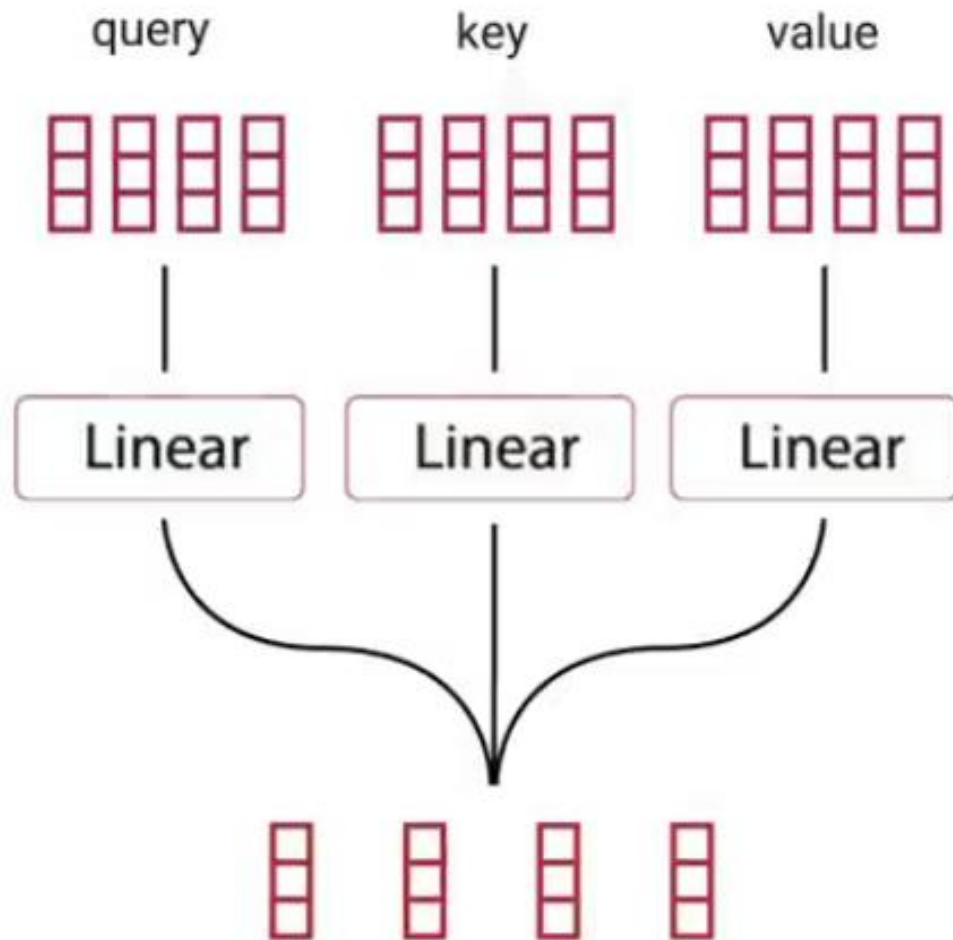4. Extract features with high attention

**Attention score:** compute pairwise similarity between each query and key

How to compute similarity between two sets of features?



Also known as the "cosine similarity"

# Learning Self-Attention

$$\frac{\boxed{\phantom{xxx}}}{\sqrt{d_k}} = \boxed{\phantom{xxx}}$$

Scaled Scores

MatMul

SoftMax

Mask (opt.)

Scale

MatMul

Q    K    V

# Learning Self-Attention

$$\text{Softmax}(\boxplus) =$$

|       | Hi  | how | are | you |
|-------|-----|-----|-----|-----|
| Hi    | 0.7 | 0.1 | 0.1 | 0.1 |
| how   | 0.1 | 0.6 | 0.2 | 0.1 |
| are   | 0.1 | 0.3 | 0.6 | 0.1 |
| you   | 0.1 | 0.3 | 0.3 | 0.3 |

MatMul

SoftMax

Mask (opt.)

Scale

MatMul

Q   K   V

# Self-Attention

Attention : What part of the input should we focus?

Attention Vectors

Focus

The → The big red dog $[0.71 \quad 0.04 \quad 0.07 \quad 0.18]^T$

big → The big red dog $[0.01 \quad 0.84 \quad 0.02 \quad 0.13]^T$

red → The big red dog $[0.09 \quad 0.05 \quad 0.62 \quad 0.24]^T$

dog → The big red dog $[0.03 \quad 0.03 \quad 0.03 \quad 0.91]^T$

# Learning Self-Attention

**Goal:** identify and attend to most important features in input.

1. Encode **position** information
2. Extract **query, key, value** for search
3. Compute **attention weighting**
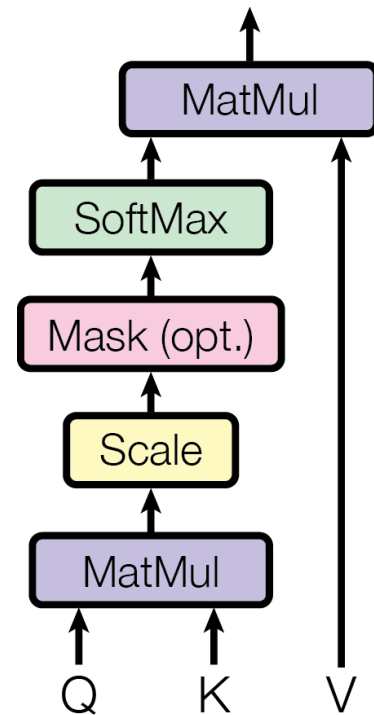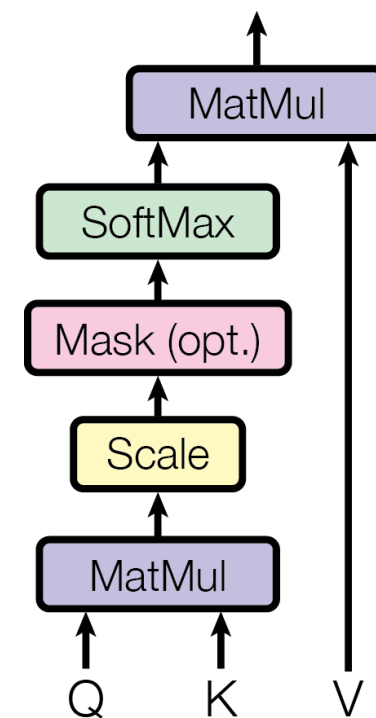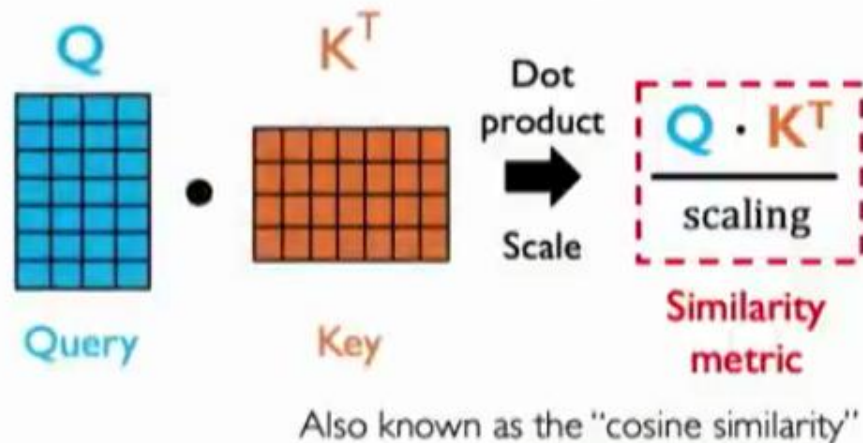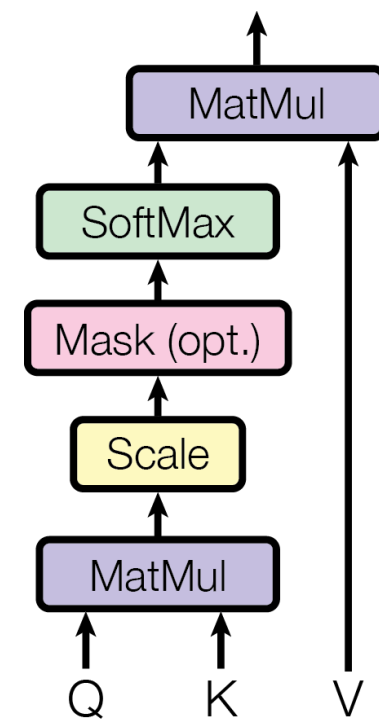4. Extract **features with high attention**

**Last step:** self-attend to extract features



Attention weighting    Value    Output

$$softmax\left(\frac{Q \cdot K^T}{scaling}\right) \cdot V = A(Q, K, V)$$



MatMul

SoftMax

Mask (opt.)

Scale

MatMul

Q    K    V

# Self-Attention



$$Z = softmax\left(\frac{Q.K^T}{\sqrt{Dimension\ of\ vector\ Q, K\ or\ V}}\right).V$$

# Multi-Head Self-Attention



Attention weighting × Value = Output

Output of attention head 1

Output of attention head 2

Output of attention head 3

# Multi-Head Self-Attention



$$Z = softmax\left(\frac{Q.K^T}{\sqrt{Dimension\ of\ vector\ Q, K\ or\ V}}\right).V$$

# Multi-Head Self-Attention



$$Z = softmax\left(\frac{Q.K^T}{\sqrt{Dimension\ of\ vector\ Q, K\ or\ V}}\right).V$$

# Residual, Add, Normalization

# Residual Layer

By explicitly learning the residual mapping, the network can focus on learning the fine-grained details or changes needed to refine the input, rather than trying to learn the complete transformation from scratch.

# Layer Normalization



1 Batch with 3 samples

| Features | | | |
|---|---|---|---|
| $x\_1$ | 1 | 3 | 8 |
| $x\_2$ | 3 | 4 | 3 |
| $x\_3$ | 5 | 6 | 2 |
| $x\_4$ | 7 | 2 | 1 |

mean       4      3.75   3.50
std_dev   2.23   1.47   2.69

Normalization across features,
independently for each sample

- can deal with sequences
- any batch number works
- can parallelize
- cannot work well with CNN

# Add & Layer Normalization & Feed-Forward

# Encoder

# English-French Translation

The big red dog
Le gros chien rouge

# English-French Translation

# Decoder



Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Nx

Add & Norm

Masked
Multi-Head
Attention

Nx

Positional
Encoding

Positional
Encoding

Input
Embedding

Output
Embedding

Inputs

Outputs
(shifted right)

# Masked Self-Attention

# Masked Self-Attention

# Masked Self-Attention

# Masked Self-Attention

# Masked Self-Attention

# Masked Self-Attention

Decoder

Self
Attention

Le $\rightarrow$ Le gros chien rouge

gros $\rightarrow$ Le gros chien rouge

chien $\rightarrow$ Le gros chien rouge

rouge $\rightarrow$ Le gros chien rouge

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 0.1 \\ 0.9 \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 0.05 \\ 0.40 \\ 0.55 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 0.16 \\ 0.09 \\ 0.15 \\ 0.66 \end{bmatrix}$$

# Encoder-Decoder Attention

# Encoder-Decoder Attention

<u>Decoder</u>

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$
Le

$$\begin{bmatrix} 0.1 \\ 0.9 \\ 0 \\ 0 \end{bmatrix}$$
gros

$$\begin{bmatrix} 0.05 \\ 0.40 \\ 0.55 \\ 0 \end{bmatrix}$$
chien

$$\begin{bmatrix} 0.16 \\ 0.09 \\ 0.15 \\ 0.66 \end{bmatrix}$$
rouge

Encapsulates English-French Interactions

Encoder-Decoder Attention

$$\begin{bmatrix} 0.71 \\ 0.04 \\ 0.07 \\ 0.18 \end{bmatrix}$$
The

$$\begin{bmatrix} 0.01 \\ 0.84 \\ 0.02 \\ 0.13 \end{bmatrix}$$
big

$$\begin{bmatrix} 0.09 \\ 0.05 \\ 0.62 \\ 0.24 \end{bmatrix}$$
red

$$\begin{bmatrix} 0.03 \\ 0.03 \\ 0.03 \\ 0.91 \end{bmatrix}$$
dog

<u>Encoder</u>

# Encoder-Decoder Attention

The memory **keys** and **values** come from the output of the **encoder**.

The **queries** come from the previous **decoder** layer.

# Decoder



Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Nx

Add & Norm

Masked
Multi-Head
Attention

Positional
Encoding

Output
Embedding

Positional
Encoding

Outputs
(shifted right)

Add & Norm

Feed
Forward

Nx

Add & Norm

Multi-Head
Attention

Positional
Encoding

Input
Embedding

Inputs

# Decoder



Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

N×

N×

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

# Decoder



| 0.1 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ...... | 0.0 | 0.0 | 0.8 |   N Class (vocab size)

**Softmax**

N Class (vocab size)

**Linear (classifier)**

Decoder

The big red dog
Le gros chien rouge

Outputs (shifted right) → Output Embedding → ⊕ ← Positional Encoding

Masked Multi-Head Attention → Add & Norm → Multi-Head Attention → Add & Norm → Feed Forward → Add & Norm

Nx

Linear → Softmax → Output Probabilities

Probability Distribution

Le → gros

Next Predicted Word

Feed Forward Layer
# Neurons =
# words in French

# Decoder

The big red dog
Le gros chien rouge



Probability Distribution

chien — Outputs (shifted right) → Output Embedding ⊕ Positional Encoding → [ Masked Multi-Head Attention → Add & Norm → Multi-Head Attention → Add & Norm → Feed Forward → Add & Norm ] Nx → Linear → Softmax → Output Probabilities → rouge

Next Predicted Word

Feed Forward Layer
# Neurons =
# words in French

Transformers
Encoder

Transformers
Decoder

# Results

Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

| Model | BLEU | | Training Cost (FLOPs) | |
|---|---|---|---|---|
| | EN-DE | EN-FR | EN-DE | EN-FR |
| ByteNet [15] | 23.75 | | | |
| Deep-Att + PosUnk [32] | | 39.2 | | $1.0 \cdot 10^{20}$ |
| GNMT + RL [31] | 24.6 | 39.92 | $2.3 \cdot 10^{19}$ | $1.4 \cdot 10^{20}$ |
| ConvS2S [8] | 25.16 | 40.46 | $9.6 \cdot 10^{18}$ | $1.5 \cdot 10^{20}$ |
| MoE [26] | 26.03 | 40.56 | $2.0 \cdot 10^{19}$ | $1.2 \cdot 10^{20}$ |
| Deep-Att + PosUnk Ensemble [32] | | 40.4 | | $8.0 \cdot 10^{20}$ |
| GNMT + RL Ensemble [31] | 26.30 | 41.16 | $1.8 \cdot 10^{20}$ | $1.1 \cdot 10^{21}$ |
| ConvS2S Ensemble [8] | 26.36 | **41.29** | $7.7 \cdot 10^{19}$ | $1.2 \cdot 10^{21}$ |
| Transformer (base model) | 27.3 | 38.1 | $\mathbf{3.3 \cdot 10^{18}}$ | |
| Transformer (big) | **28.4** | **41.0** | $2.3 \cdot 10^{19}$ | |

# Thank you



Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Add & Norm

Masked
Multi-Head
Attention

N×

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

N×

Positional
Encoding

Positional
Encoding

Input
Embedding

Output
Embedding

Inputs

Outputs
(shifted right)