# 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
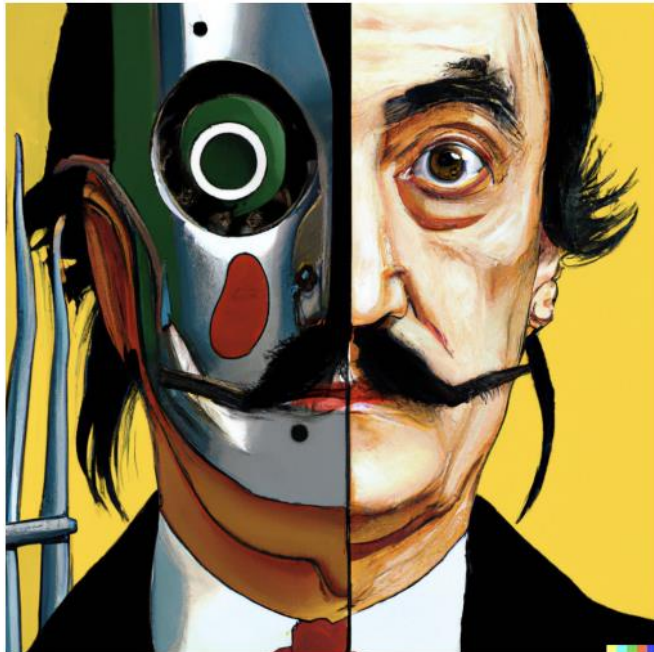aidan@cs.toronto.edu

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

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

## Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention

m



vibrant portrait painting of Salvador Dalí with a robotic half face

a shiba inu wearing a beret and black turtleneck

a close up of a handpalm with leaves growing from it

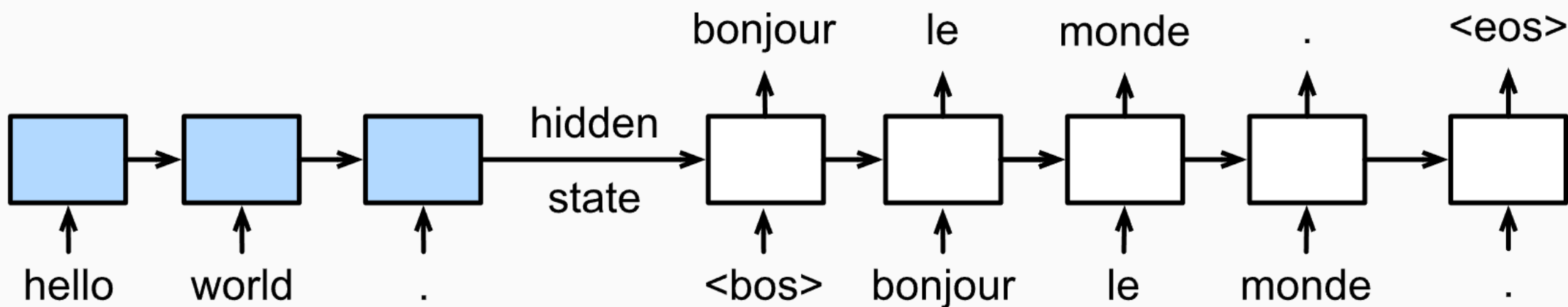an espresso machine that makes coffee from human souls, artstation

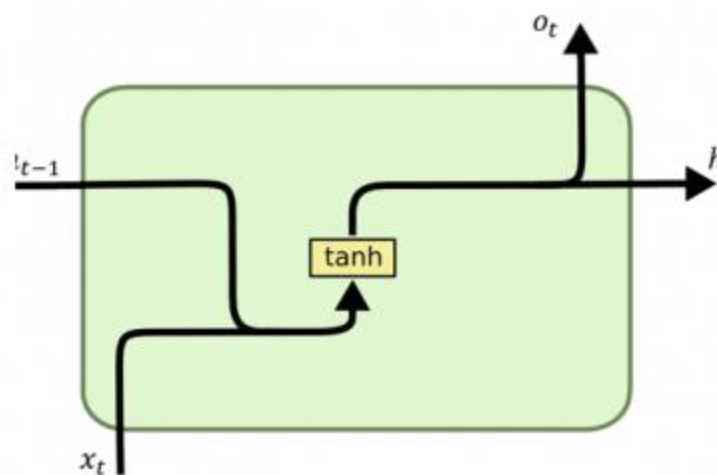panda mad scientist mixing sparkling chemicals, artstation

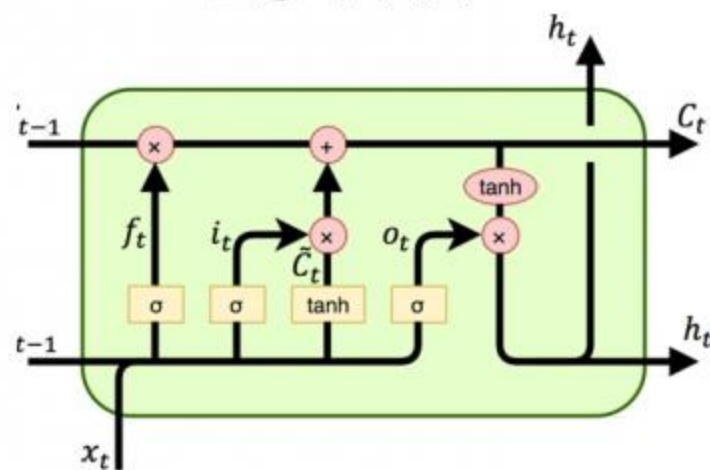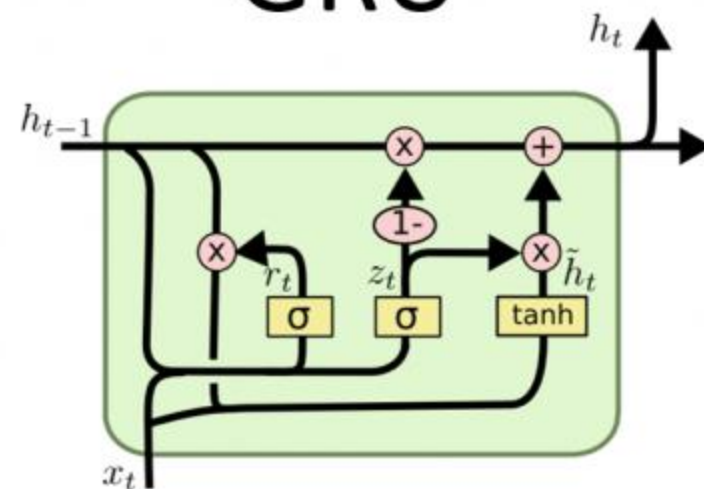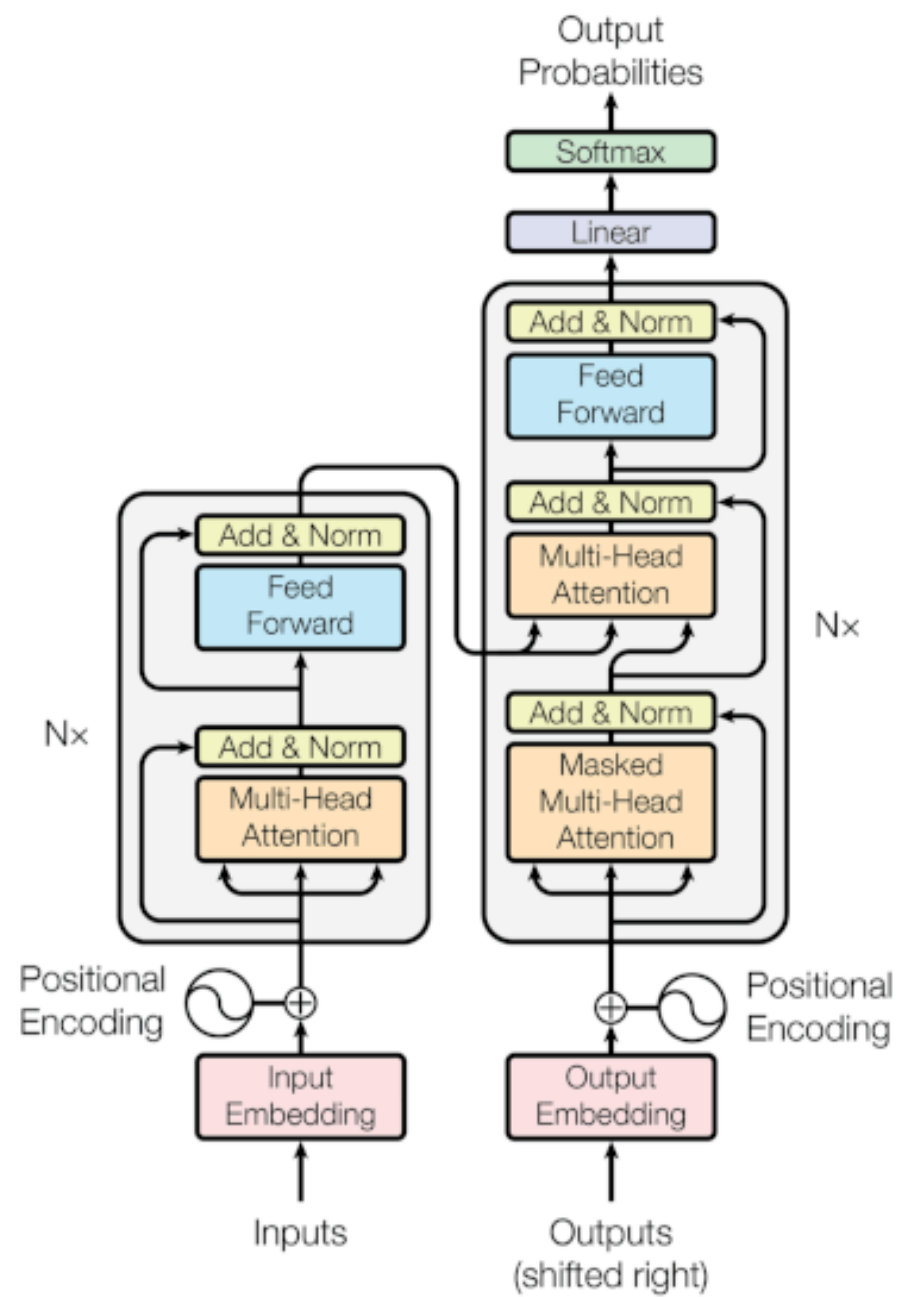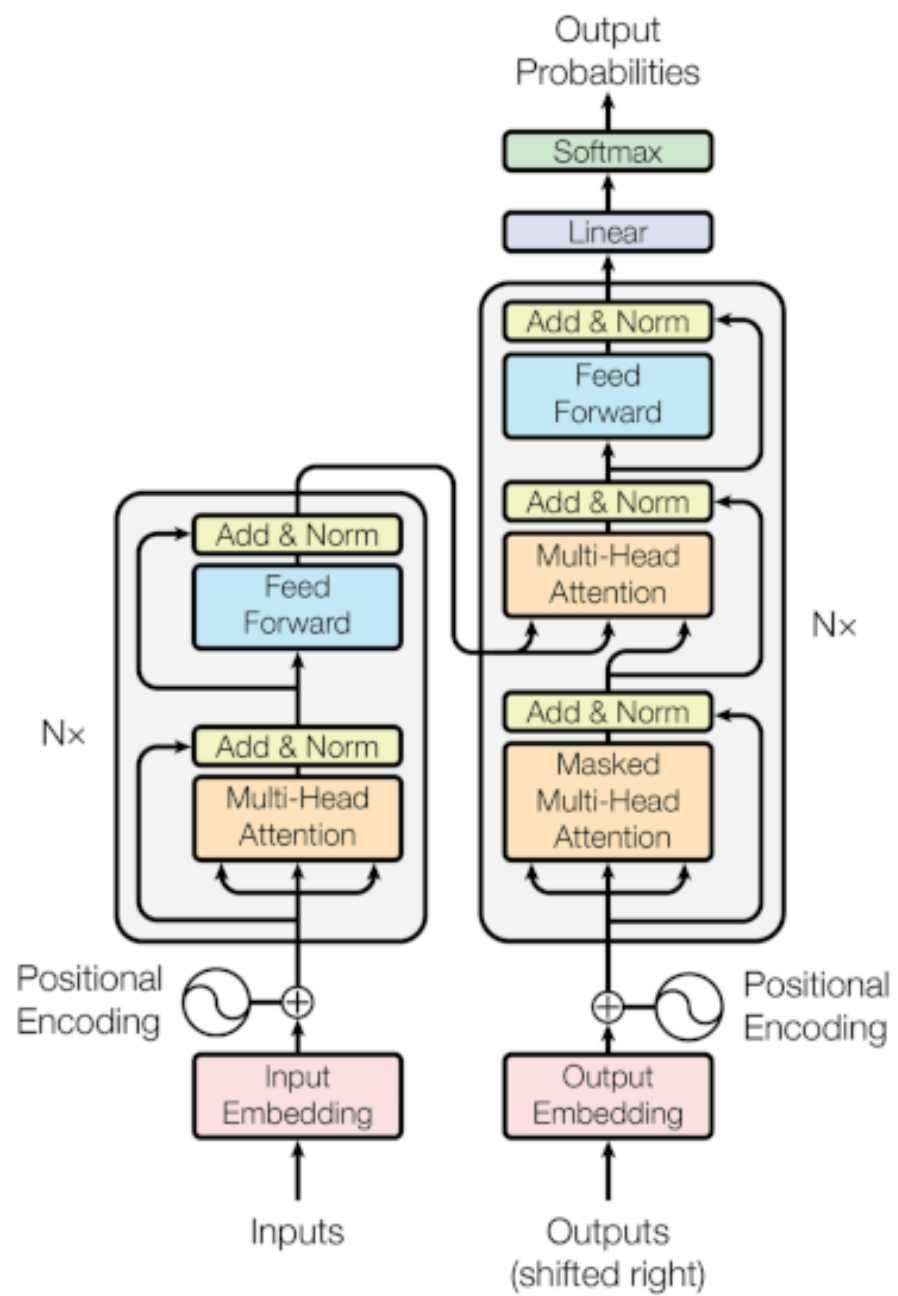a corgi's head depicted as an explosion of a nebula

Encoder ... hidden state ... Decoder

Encoder inputs: hello, world, .

Decoder outputs: bonjour, le, monde, ., <eos>

Decoder inputs: <bos>, bonjour, le, monde, .

# RNN

$o_t$

$h_{t-1}$

tanh

$x_t$

# LSTM

$h_t$

$C_{t-1}$

$f_t$ $i_t$ $\tilde{C}_t$ $o_t$

$\times$ $+$ tanh $\times$

$\sigma$ $\sigma$ tanh $\sigma$

$C_t$

$h_t$

$x_t$

# GRU

$h_t$

$h_{t-1}$

$r_t$ $z_t$ $\tilde{h}_t$

$\times$ $1-$ $\times$ $+$

$\sigma$ $\sigma$ tanh

$x_t$

Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Add & Norm

Masked
Multi-Head
Attention

Nx

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Nx

Positional
Encoding

Positional
Encoding

Input
Embedding

Output
Embedding

Inputs

Outputs
(shifted right)

INPUT

Je    suis    étudiant

THE
TRANSFORMER

OUTPUT

I    am    a    student

**Output Probabilities**

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

N×

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

OUTPUT: I am a student

ENCODERS → DECODERS

INPUT: Je suis étudiant

Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Nx

Add & Norm

Masked
Multi-Head
Attention

Add & Norm

Feed
Forward

Nx

Add & Norm

Multi-Head
Attention

Positional
Encoding

Positional
Encoding

Input
Embedding

Output
Embedding

Inputs

Outputs
(shifted right)

OUTPUT    I   am   a   student

ENCODERS → DECODERS

INPUT    Je   suis   étudiant

Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Nx

Add & Norm

Feed
Forward

Add & Norm

Add & Norm

Multi-Head
Attention

Masked
Multi-Head
Attention

Nx

Positional
Encoding

Positional
Encoding

Input
Embedding

Output
Embedding

Inputs

Outputs
(shifted right)

OUTPUT    I   am   a   student

ENCODERS

DECODERS

INPUT    Je   suis   étudiant

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Nx

Add & Norm

Feed Forward

Nx

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

OUTPUT    I    am    a    student

ENCODERS    DECODERS

INPUT    Je    suis    étudiant

Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Nx

Add & Norm

Masked
Multi-Head
Attention

Add & Norm

Feed
Forward

Nx

Add & Norm

Multi-Head
Attention

Positional
Encoding

Positional
Encoding

Input
Embedding

Output
Embedding

Inputs

Outputs
(shifted right)

OUTPUT    I    am    a    student

ENCODERS    →    DECODERS

INPUT    Je    suis    étudiant

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Nx

Add & Norm

Masked Multi-Head Attention

Add & Norm

Feed Forward

Nx

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

Positional Encoding

Output Embedding

Outputs (shifted right)

OUTPUT    I    am    a    student

ENCODERS    →    DECODERS

INPUT    Je    suis    étudiant

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Nx

Add & Norm

Feed Forward

Nx

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

OUTPUT | I | am | a | student

ENCODER → DECODER

ENCODER → DECODER

ENCODER → DECODER

ENCODER → DECODER

ENCODER → DECODER

ENCODER → DECODER

INPUT | Je | suis | étudiant

Decoding time step: 1 (2) 3 4 5 6     OUTPUT     I

K encdec     V encdec     Linear + Softmax

ENCODERS     DECODERS

EMBEDDING WITH TIME SIGNAL

EMBEDDINGS

INPUT     Je     suis     étudiant          PREVIOUS OUTPUTS     I

Je     suis     etudiant

Je                              suis                              etudiant
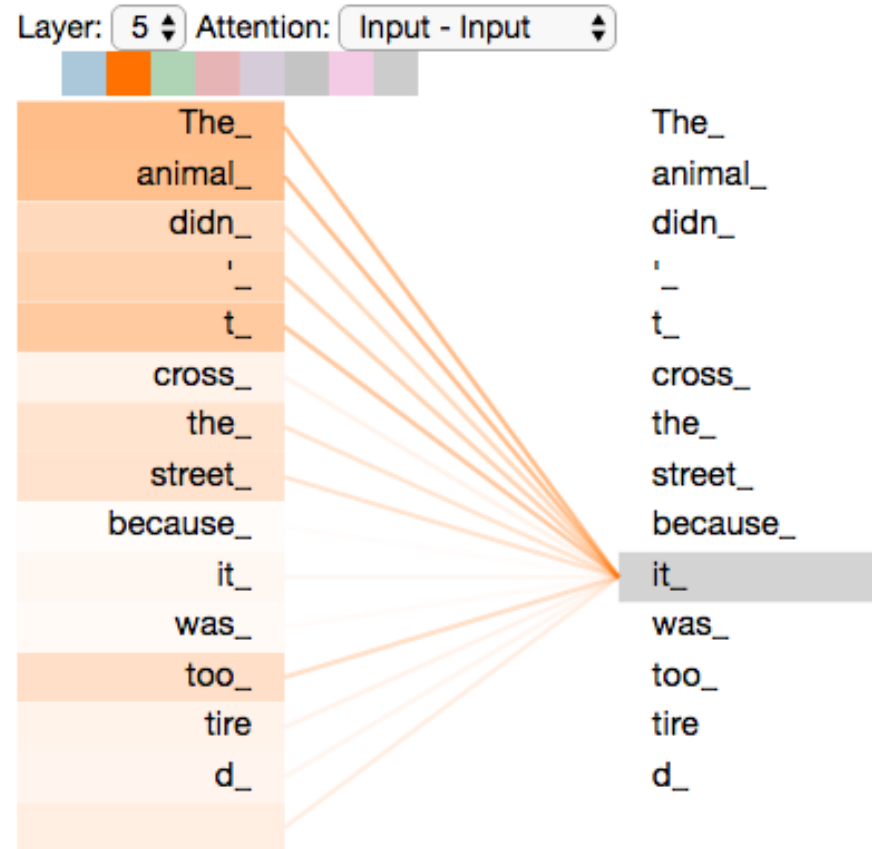
Je                                        suis                                    etudiant

Je                              suis                          etudiant

Je                                    suis                                  etudiant

Je                    suis                    etudiant

Je                    suis                    etudiant

Je                    suis                    etudiant

movie $m$

has romance
has action
has comedy

user $u$

likes romance
likes action
likes comedy

$\text{score} = u_1 m_1 + u_2 m_2 + u_3 m_3$

# TERMINOLOGY

Je                          suis                          etudiant

$$QK^T$$

$$K \quad Q \quad V$$

Je · · · · · · suis · · · · · · etudiant

$$softmax\left(\frac{\color{blue}{Q}\color{orange}{K^T}}{\sqrt{n}}\right)$$

K          Q          V

Je          suis          etudiant

|  | **Thinking** | **Machines** |
|---|---|---|
| Input | | |
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ |
| Divide by 8 ( $\sqrt{d_k}$ ) | 14 | 12 |
| Softmax | 0.88 | 0.12 |

Engels - gedetecteerd

Nederlands

yes we can

Ja dat kunnen we

| Engels - gedetecteerd ▾ | ⇄ | Nederlands ▾ |

| beer can | ✕ | bierblikje |

🎤 🔊

📋 🔊                    ✓ Geverifieerd
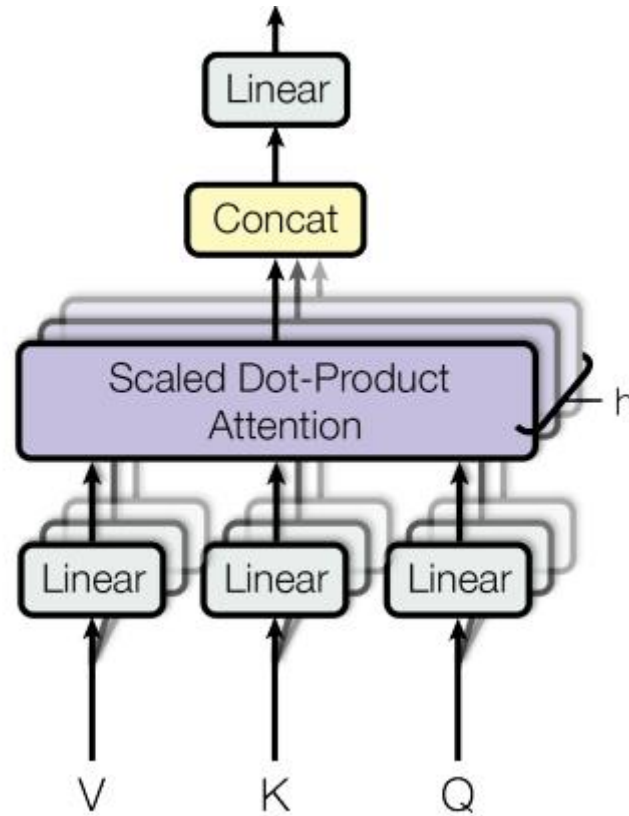
# Scaled Dot-Product Attention

Multi-Head Attention

Thinking Machines

$X$

$W_0^Q$
$W_0^K$
$W_0^V$

$Q_0$
$K_0$
$V_0$

$Z_0$

$W_1^Q$
$W_1^K$
$W_1^V$

$Q_1$
$K_1$
$V_1$

$Z_1$

...          ...          ...

$W_7^Q$
$W_7^K$
$W_7^V$

$Q_7$
$K_7$
$V_7$

$Z_7$

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

Nx

Nx

Positional
Encoding

Positional
Encoding

Input
Embedding

Output
Embedding

Inputs

Outputs
(shifted right)