



TRANSFORMER NETWORKS

Attention Is All You Need

Ashish Vaswani^{*} Noam Shazeer^{*} Niki Parmar^{*} Jakob Uszkoreit^{*}

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best

1 Introduction

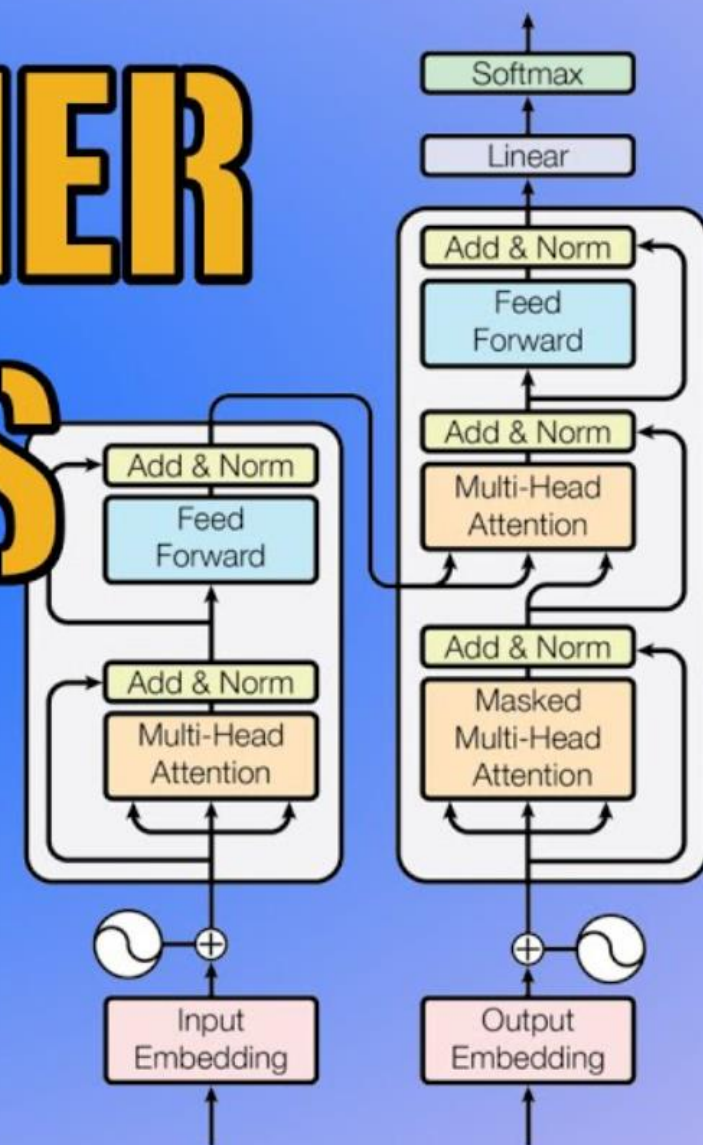
Recurrent neural networks, long short-term memory [12] and gated recurrent [7] neural networks in particular, have been firmly established as state of the art approaches in sequence modeling and transduction problems such as language modeling and machine translation [29, 2, 5]. Numerous efforts have since continued to push the boundaries of recurrent language models and encoder-decoder architectures [31, 21, 13].

^{*}Equal contribution. Listing order is random. Jakob proposed replacing RNNs with self-attention and started the effort to evaluate this idea. Ashish, with Ilia, designed and implemented the first Transformer models and has been crucially involved in every aspect of this work. Noam proposed scaled dot-product attention, multi-head attention and the parameter-free position representation and became the other person involved in nearly every detail. Niki designed, implemented, tuned and evaluated countless model variants in our original codebase and tensor2tensor. Llion also experimented with novel model variants, was responsible for our initial codebase, and efficient inference and visualizations. Lukasz and Adian spent countless long days designing various parts of and implementing tensor2tensor, replacing our earlier codebase, greatly improving results and massively accelerating our research.

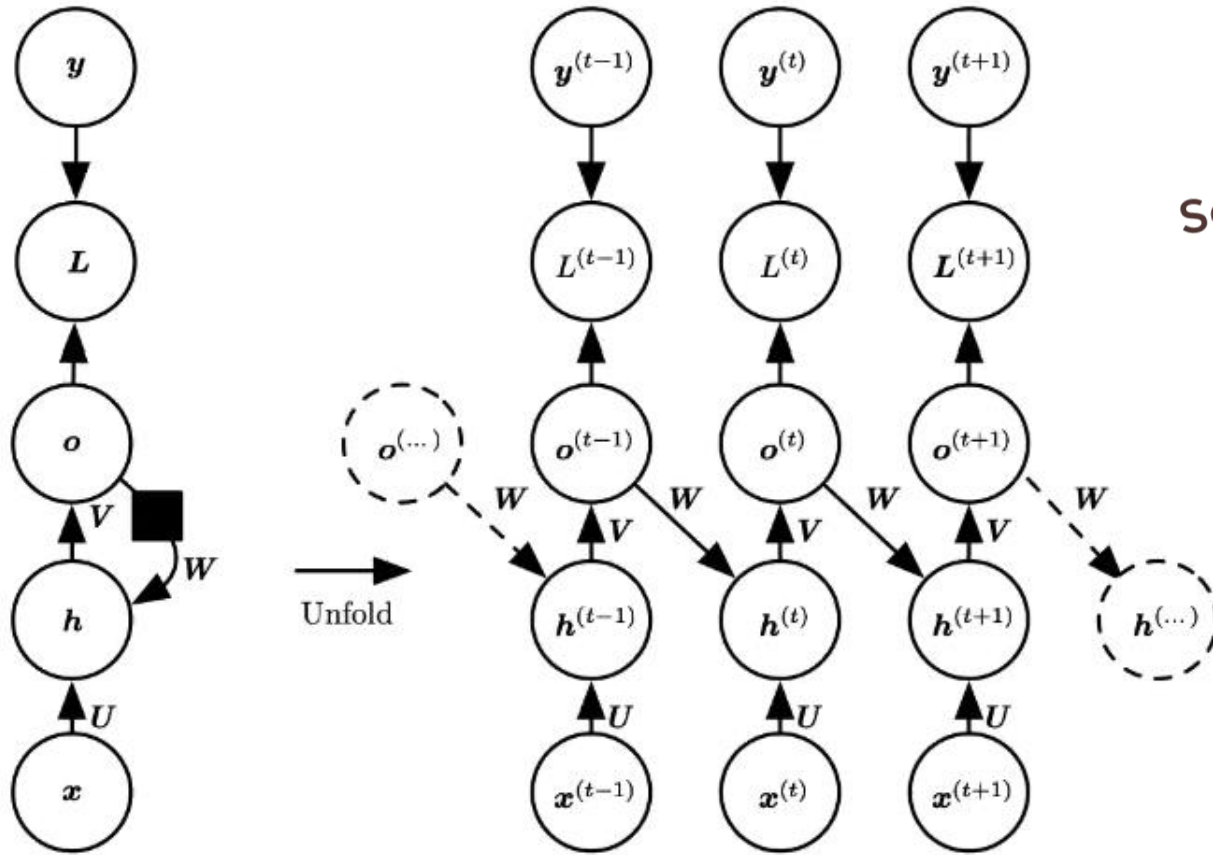
[†]Work performed while at Google Brain.

[‡]Work performed while at Google Research.

31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA.

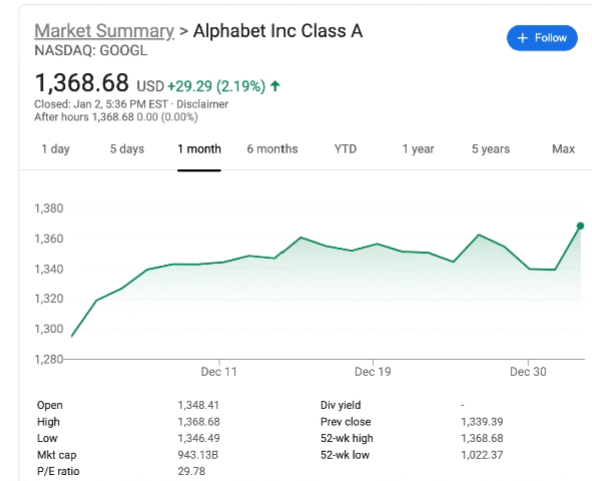


Recurrent Neural Networks



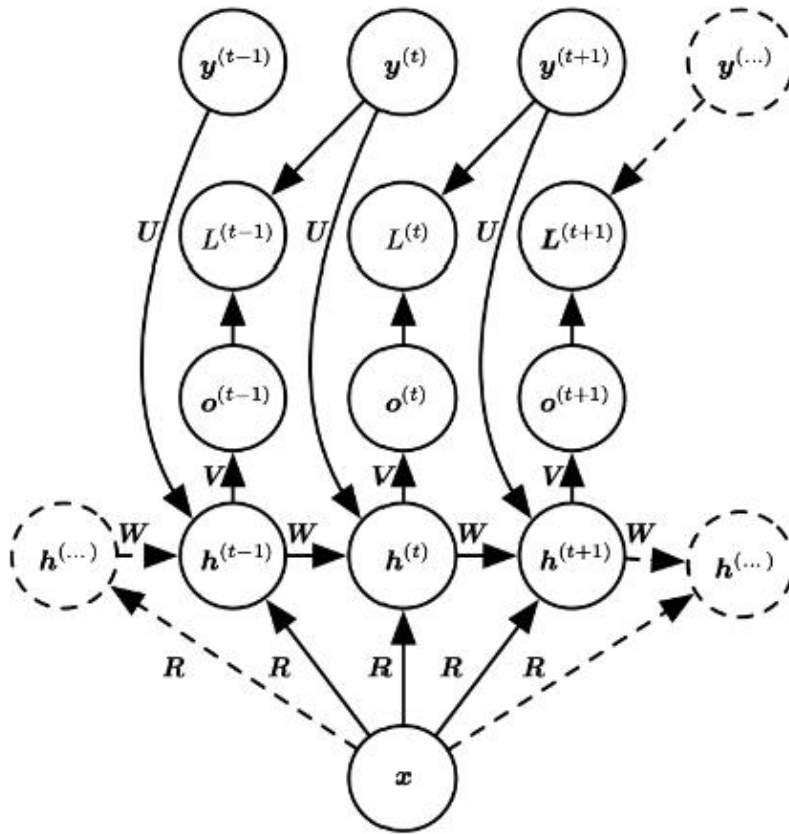
The unicorn is scotland's national animal

Sequences



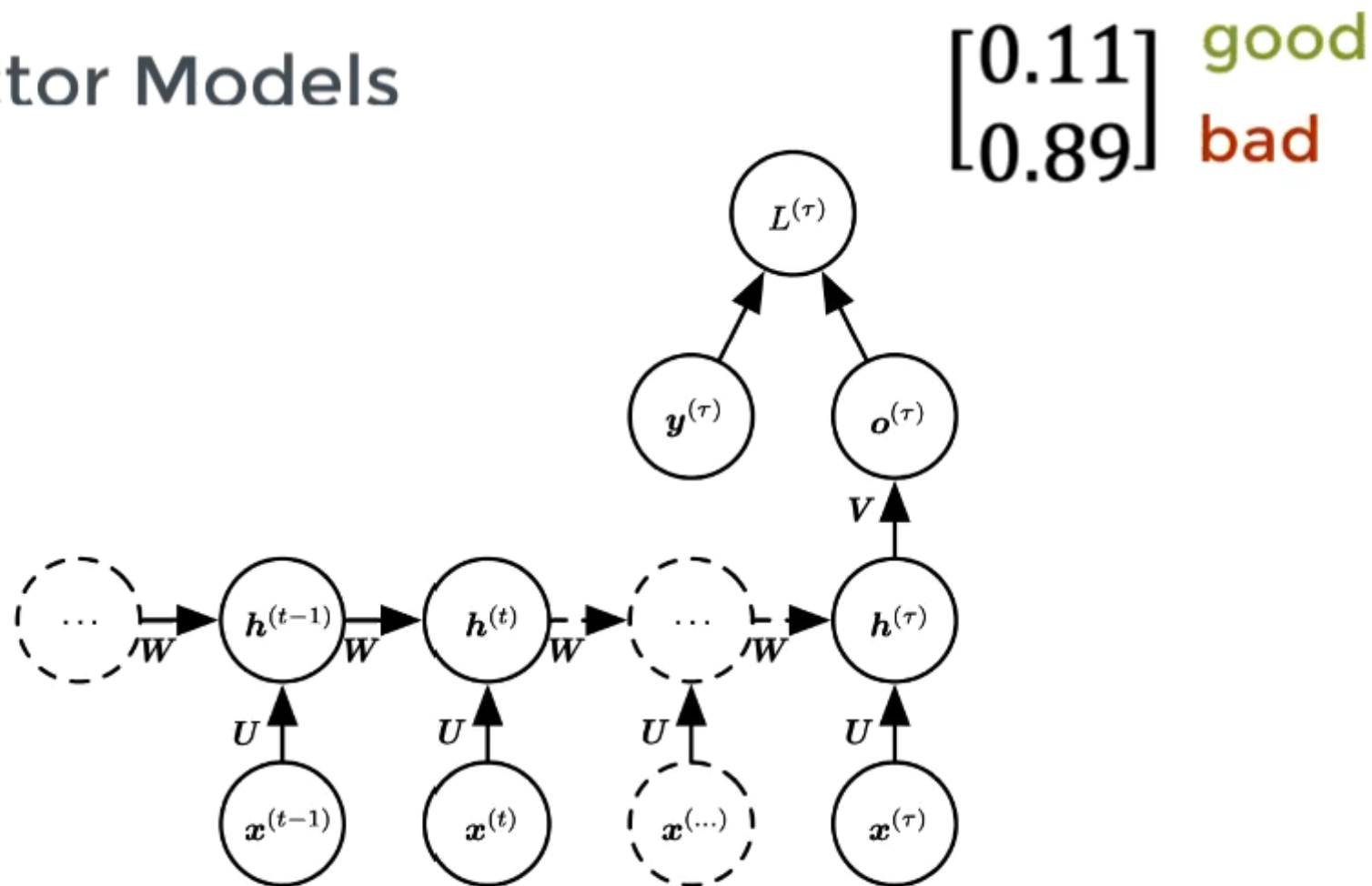
Recurrent Neural Networks

1. Vector-Sequence Models



Recurrent Neural Networks

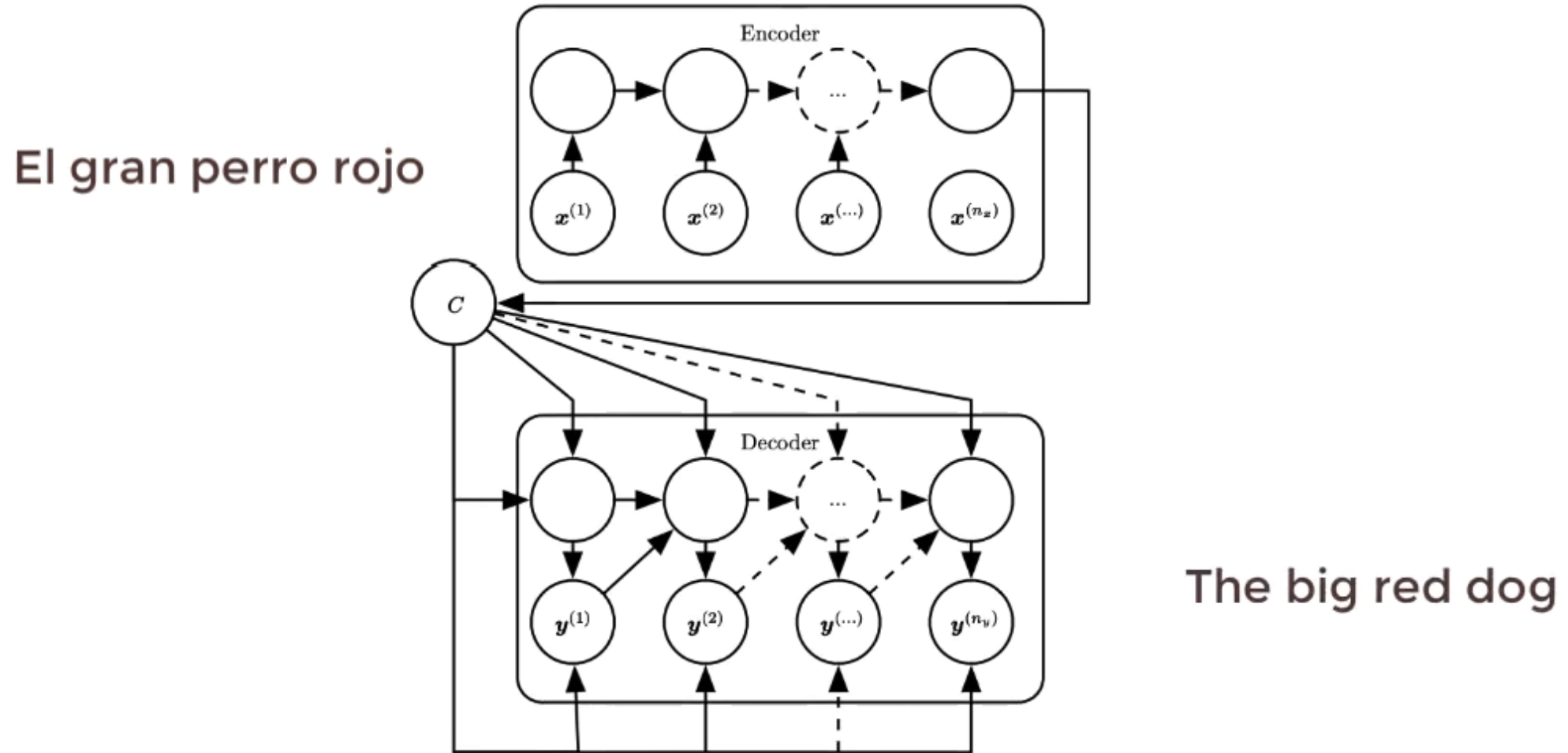
2. Sequence-Vector Models



The main character sucked

Recurrent Neural Networks

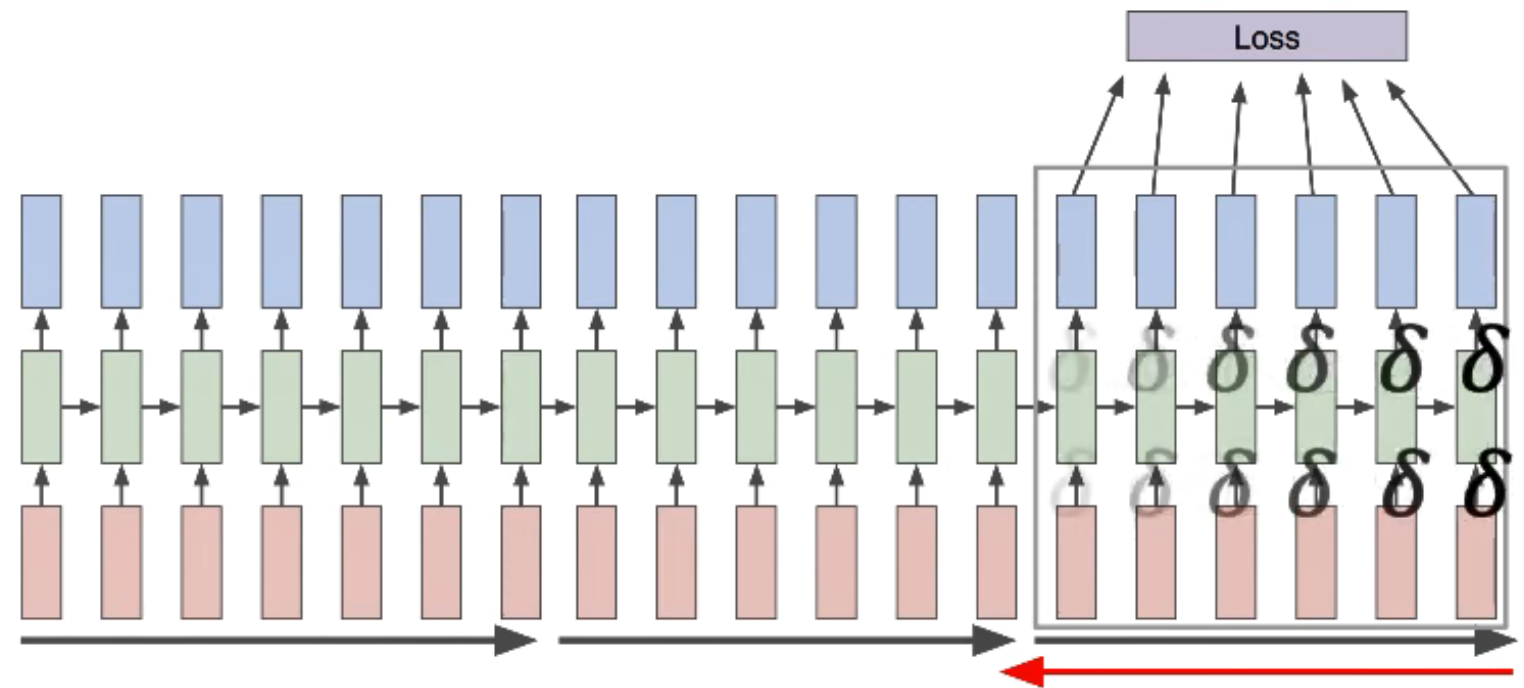
3. Sequence-Sequence Models



Recurrent Neural Networks

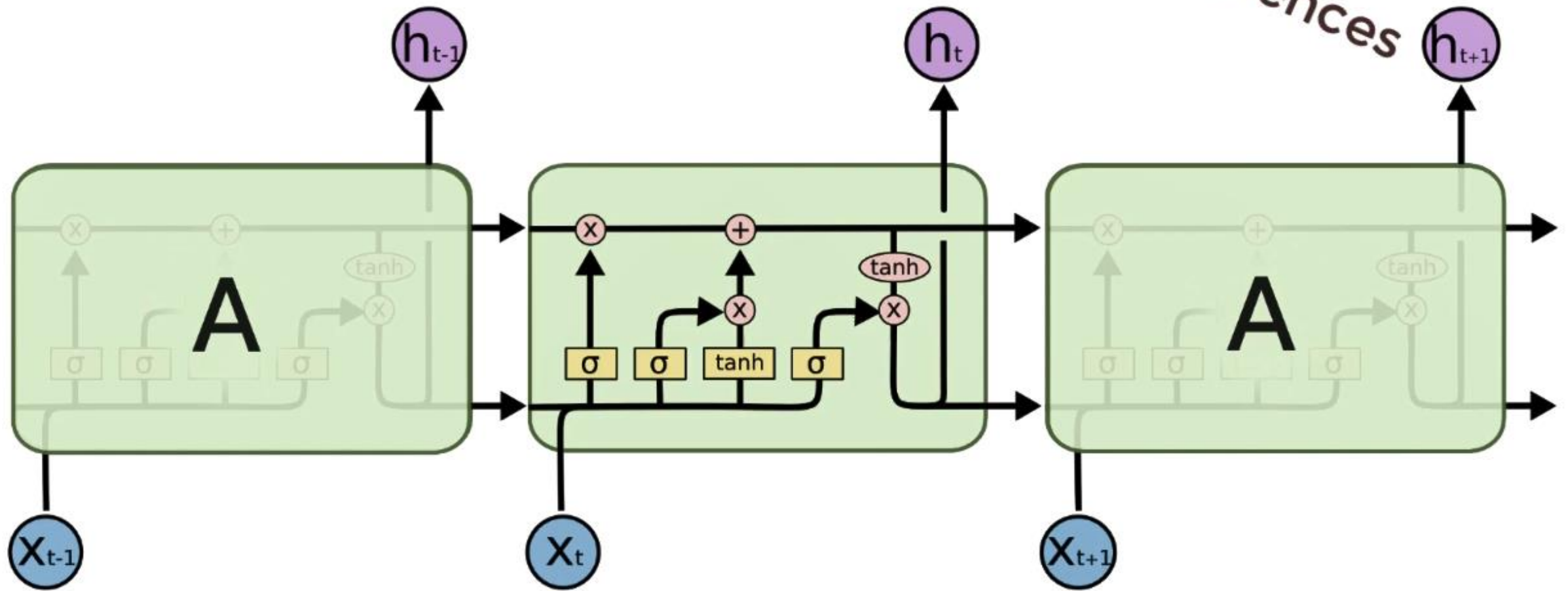
Disadvantages

1. Slow to train.
2. Long sequences lead to vanishing/exploding gradients



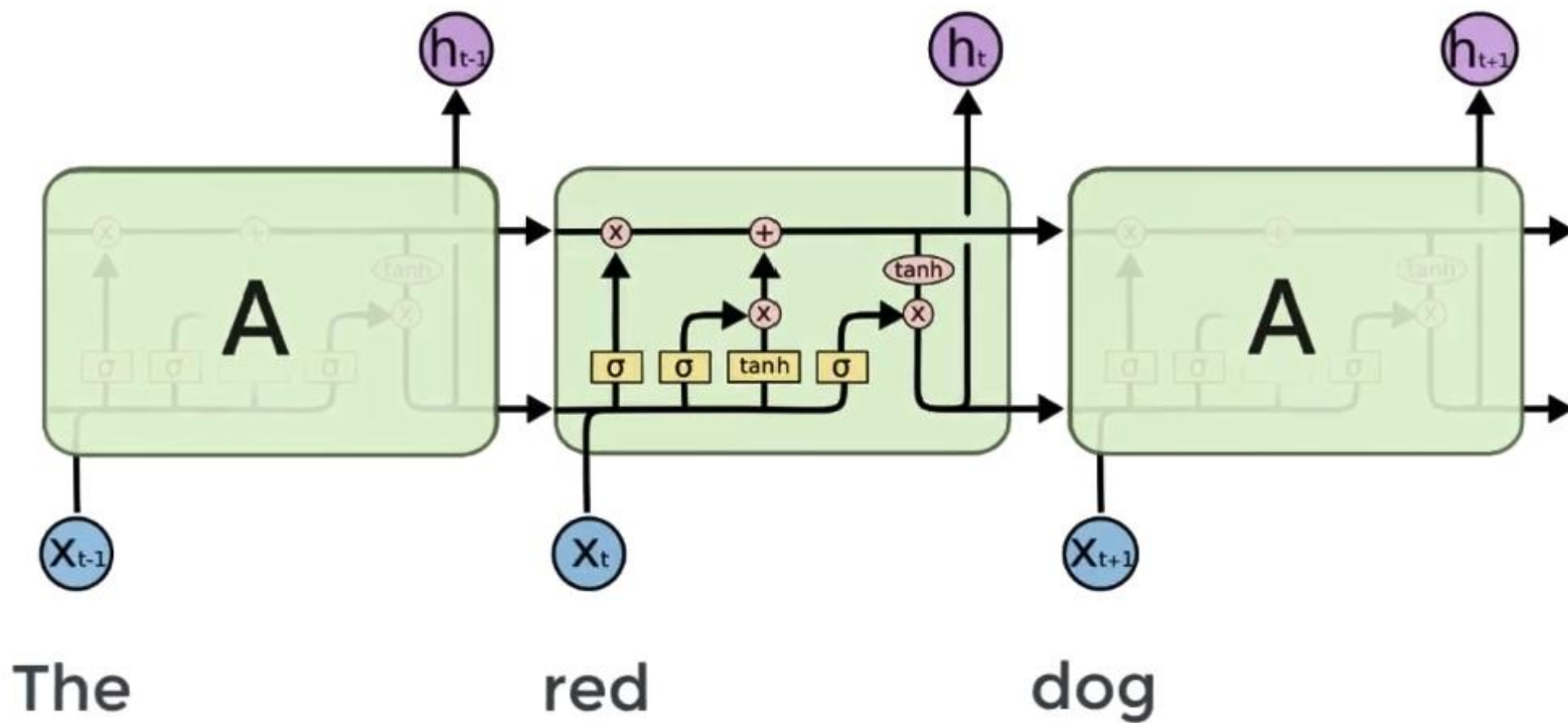
Truncated Backprop Through Time

LSTM Networks



The repeating module in an LSTM contains four interacting layers.

LSTM Networks



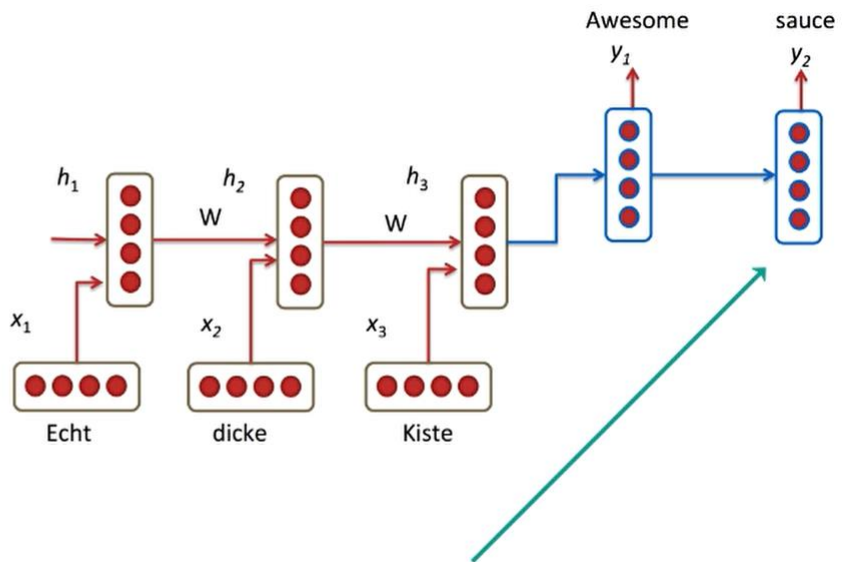
LSTM Networks

Can we parallelize
sequential data?



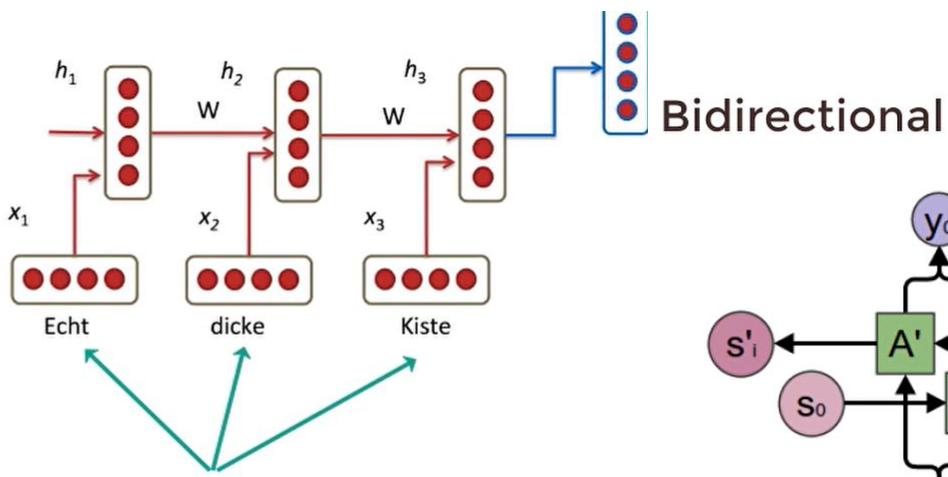
LSTM Networks

1. Slow

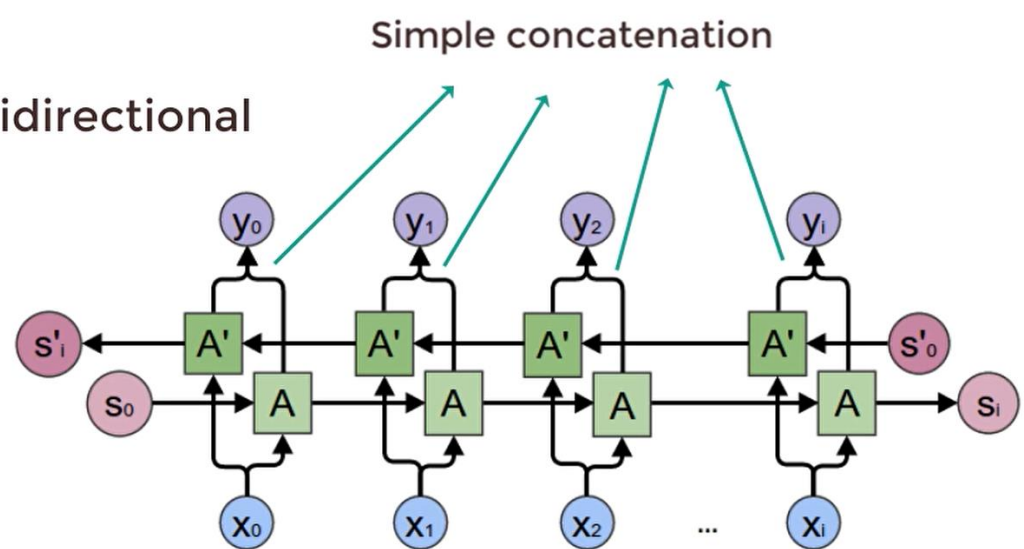


Timestep 5

2. Not truly Bidirectional

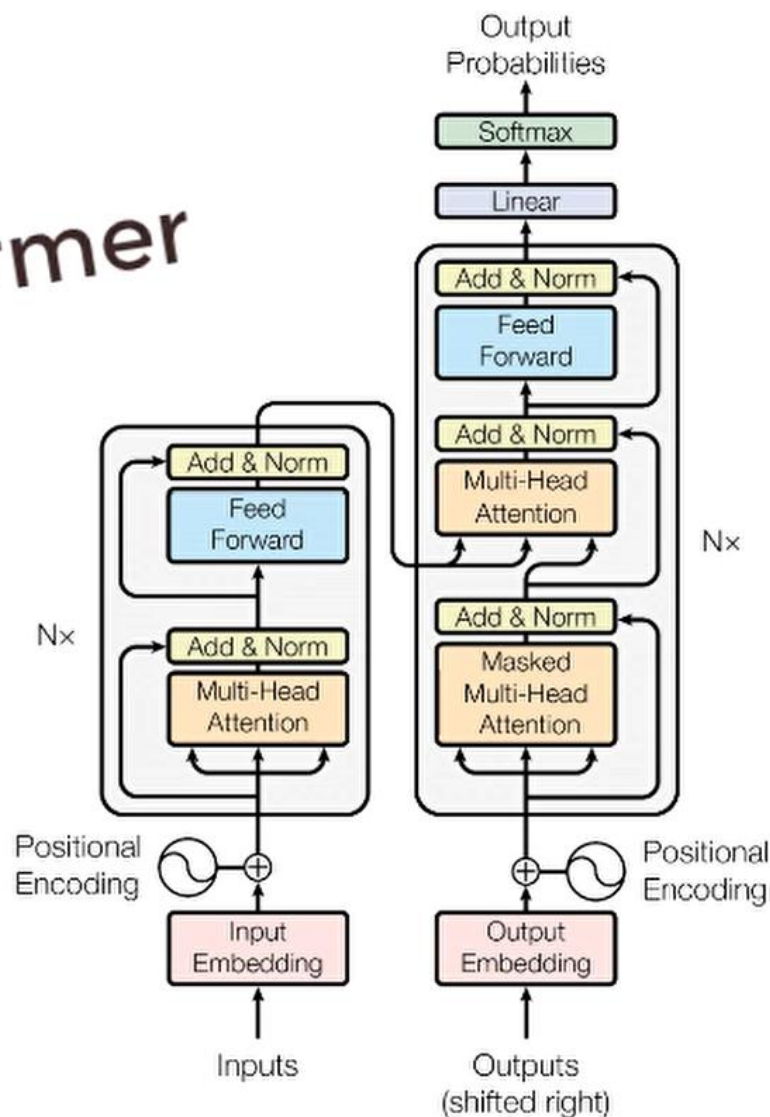


True meaning of source words not entirely captured



LSTM Vs Transformer

Transformer



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

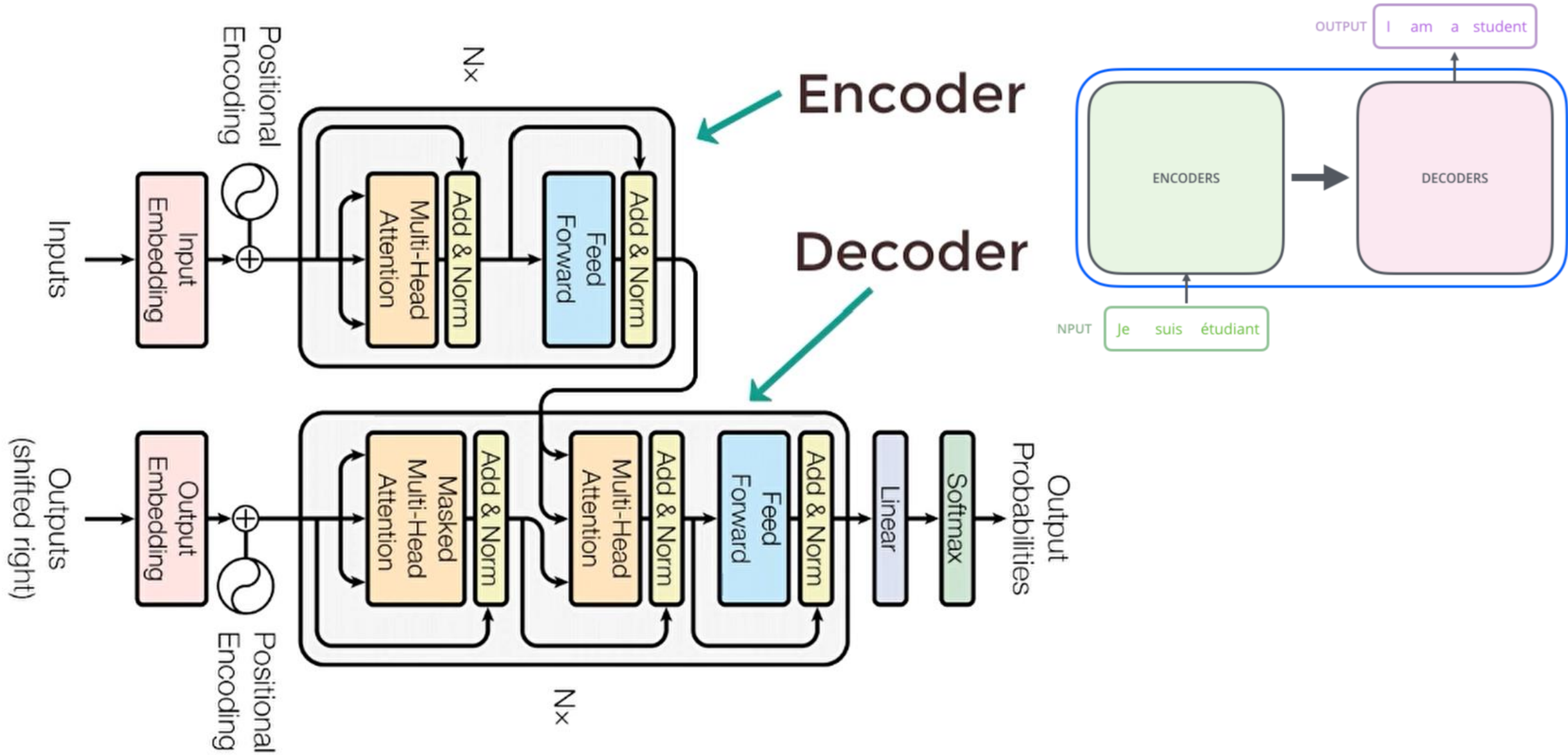
Lukasz Kaiser*
Google Brain
lukaszkaizer@google.com

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

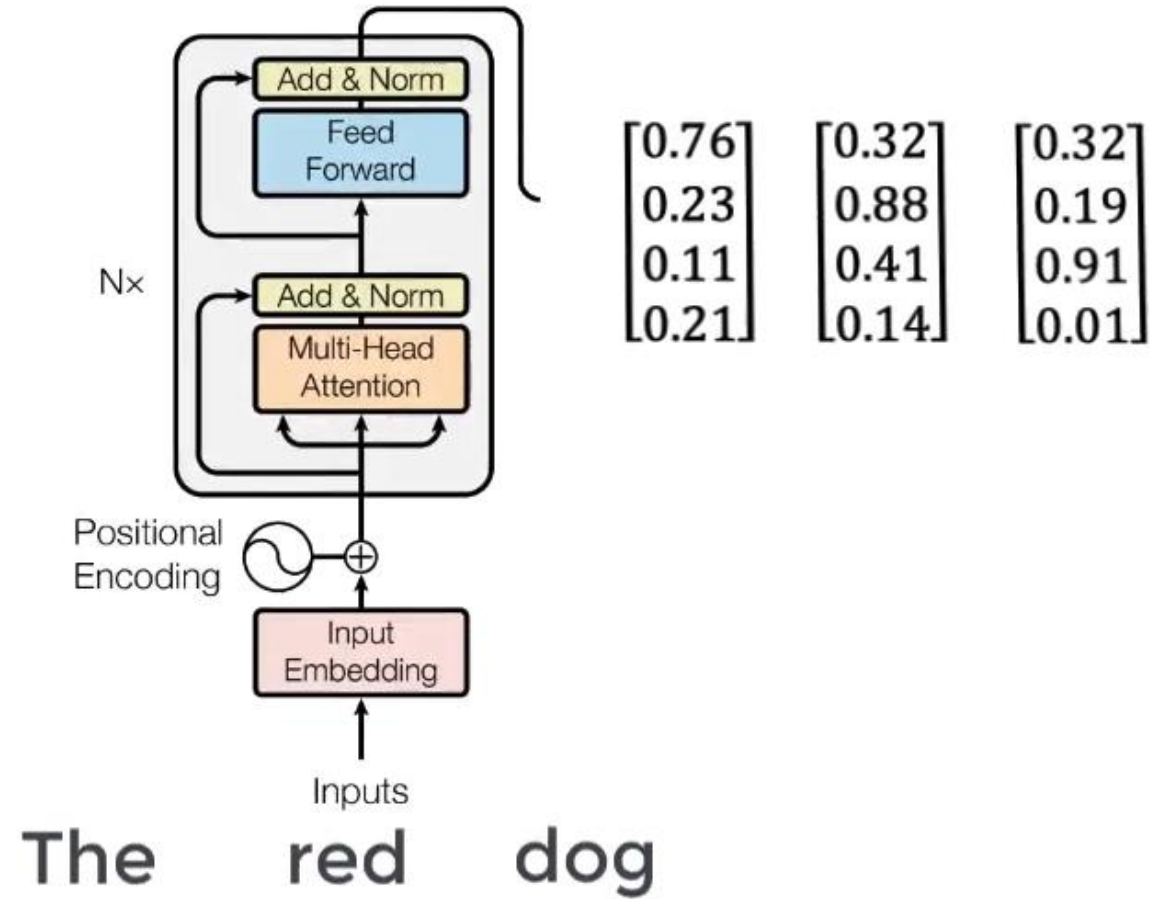
Abstract

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 mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

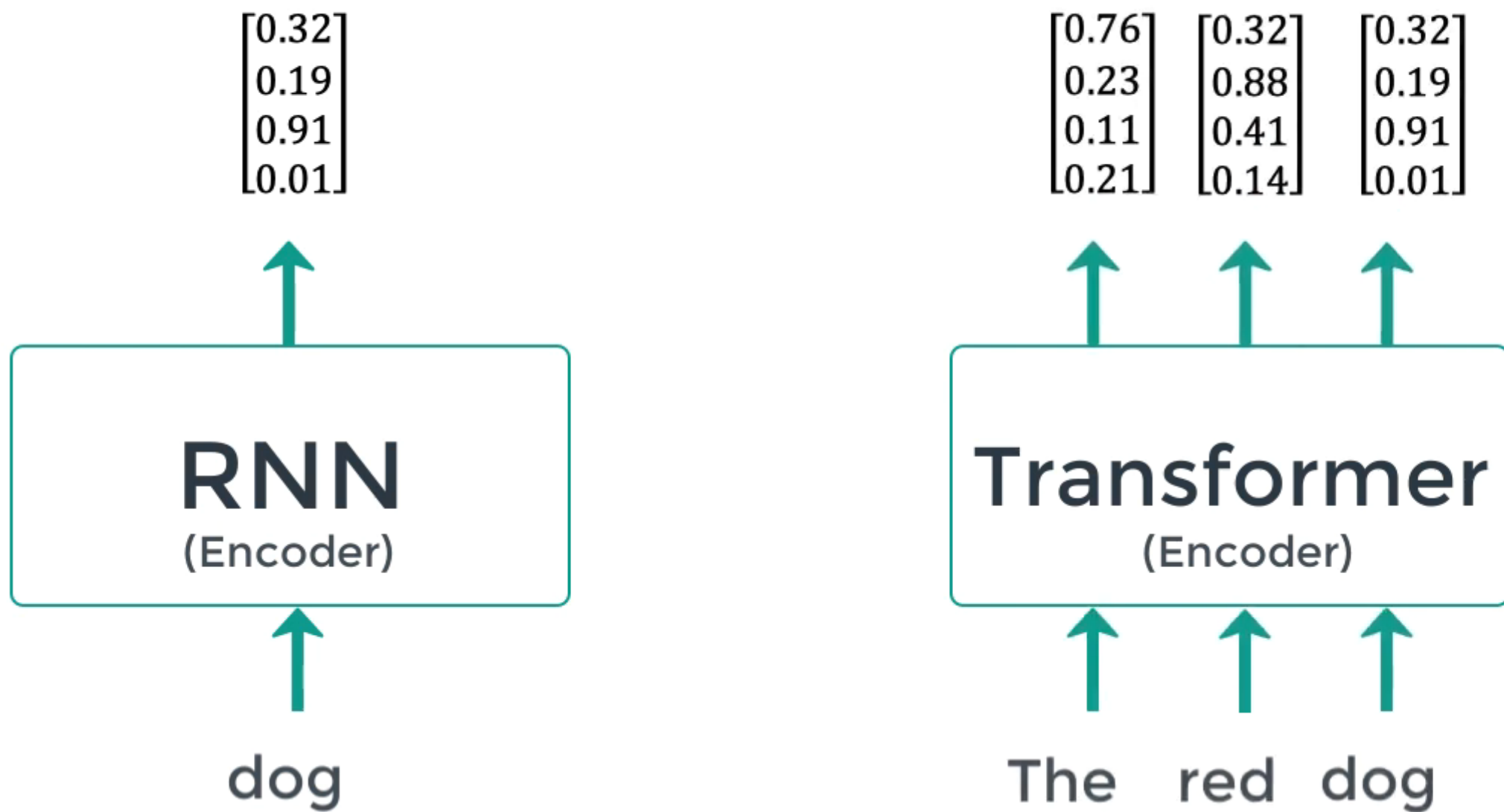
Transformer Flow



Transformers



English-French Translation



Transformer Components

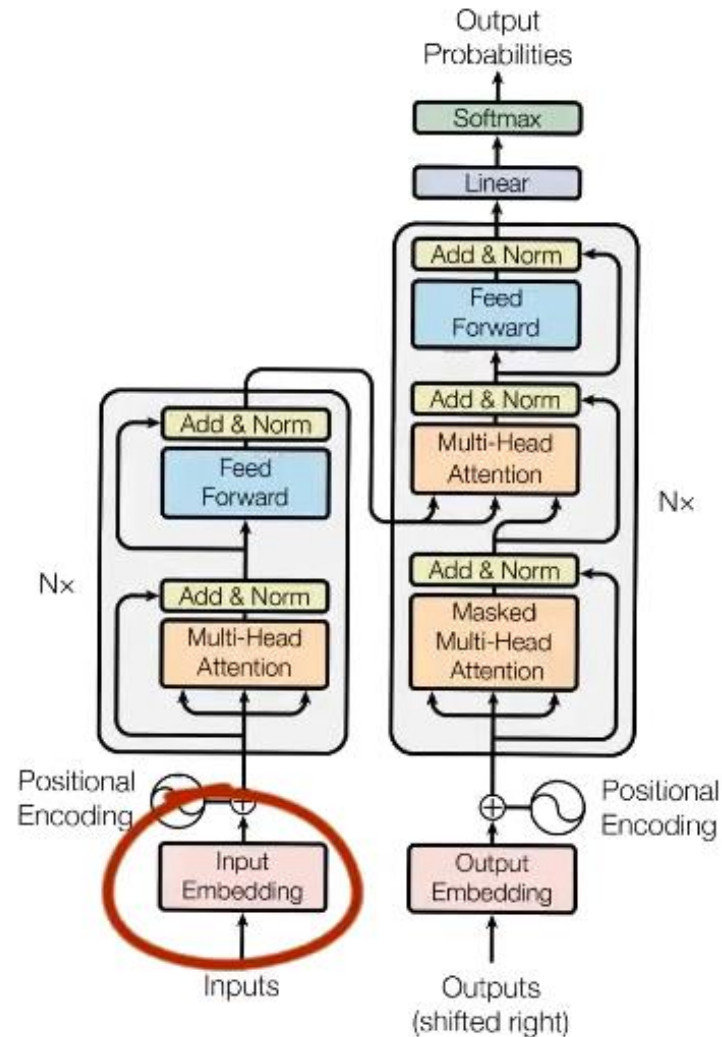


Figure 1: The Transformer - model architecture.

Transformer Components

Input Embedding





“Embedding
Space”

GloVe: Global Vectors for Word Representation

Jeffrey Pennington, Richard Socher, Christopher D. Manning

Introduction

GloVe is an unsupervised learning algorithm for obtaining vector representations for words. Training is performed on aggregated global word-word co-occurrence statistics from a corpus, and the resulting representations showcase interesting linear substructures of the word vector space.

Getting started (Code download)

- Download the [code](#) (licensed under the [Apache License, Version 2.0](#))
- Unpack the files: `unzip GloVe-1.2.zip`
- Compile the source: `cd GloVe-1.2 && make`
- Run the demo script: `./demo.sh`
- Consult the included README for further usage details, or ask a [question](#)
- The code is also available [on GitHub](#)

Download pre-trained word vectors

- Pre-trained word vectors. This data is made available under the [Public Domain Dedication and License](#) v1.0 whose full text can be found at: <http://www.opendatacommons.org/licenses/pddl/1.0/>.
 - [Wikipedia 2014](#) + [Gigaword 5](#) (6B tokens, 400K vocab, uncased, 50d, 100d, 200d, & 300d vectors, 822 MB download): [glove.6B.zip](#)
 - Common Crawl (42B tokens, 1.9M vocab, uncased, 300d vectors, 1.75 GB download): [glove.42B.300d.zip](#)
 - Common Crawl (840B tokens, 2.2M vocab, cased, 300d vectors, 2.03 GB download): [glove.840B.300d.zip](#)
 - Twitter (2B tweets, 27B tokens, 1.2M vocab, uncased, 25d, 50d, 100d, & 200d vectors, 1.42 GB download): [glove.twitter.27B.zip](#)
- Ruby [script](#) for preprocessing Twitter data

Citing GloVe

Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. [GloVe: Global Vectors for Word Representation](#). [pdf] [bib]

Highlights

1. Nearest neighbors

Transformer Components

Input Embedding

dog


$$\begin{bmatrix} 0.37 \\ 0.99 \\ 0.01 \\ 0.08 \end{bmatrix}$$

AJ's **dog** is a cutie

AJ looks like a **dog**

Transformer Components

Positional Encoder :vector that gives context based on position of word in sentence

AJ's **dog** is a cutie  Position 2

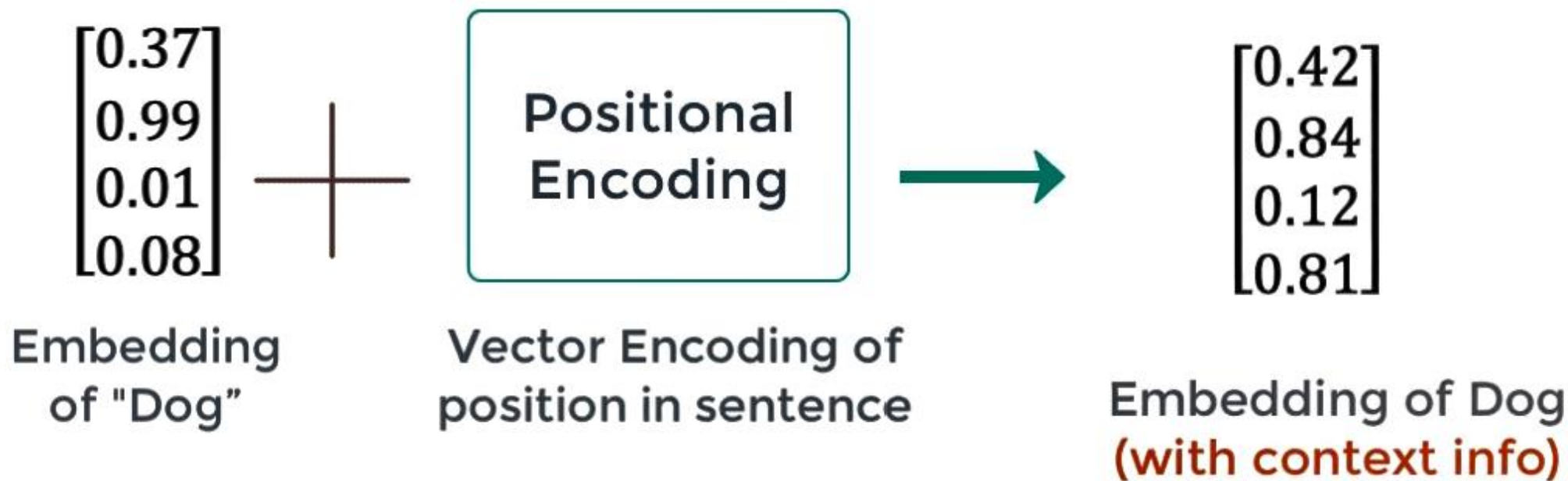
AJ looks like a **dog**  Position 5

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

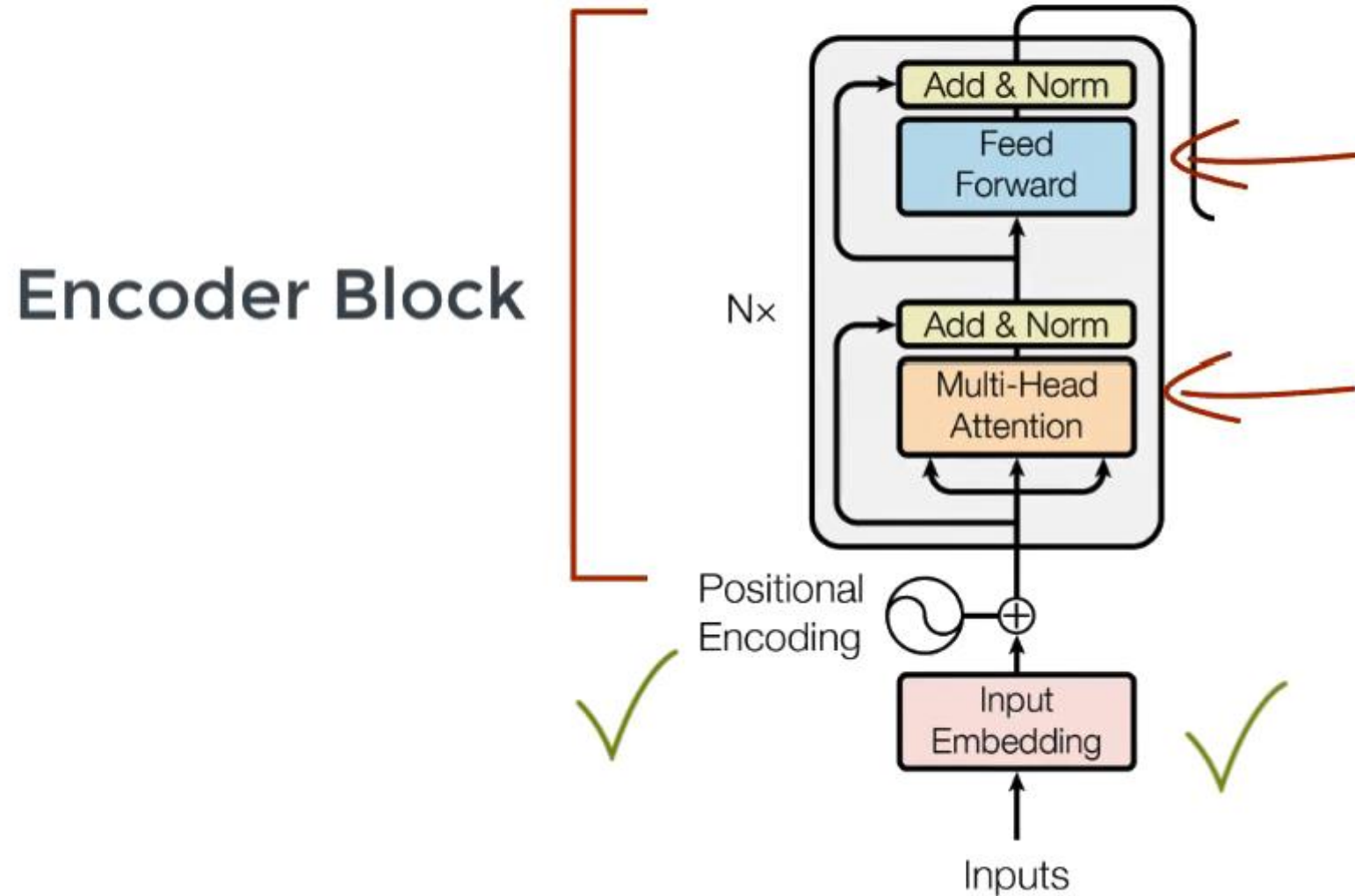
Transformer Components

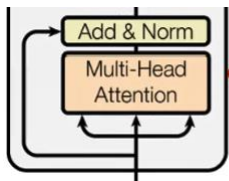
Positional Encoder :vector that gives context based on position of word in sentence



$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

Transformer Components





Transformer Components

Attention : What part of the input should we focus?

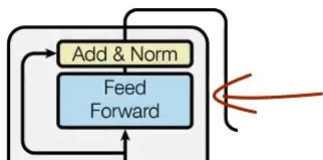
Focus

The → The big red dog
big → The big red dog
red → The big red dog
dog → The big red dog

Attention Vectors

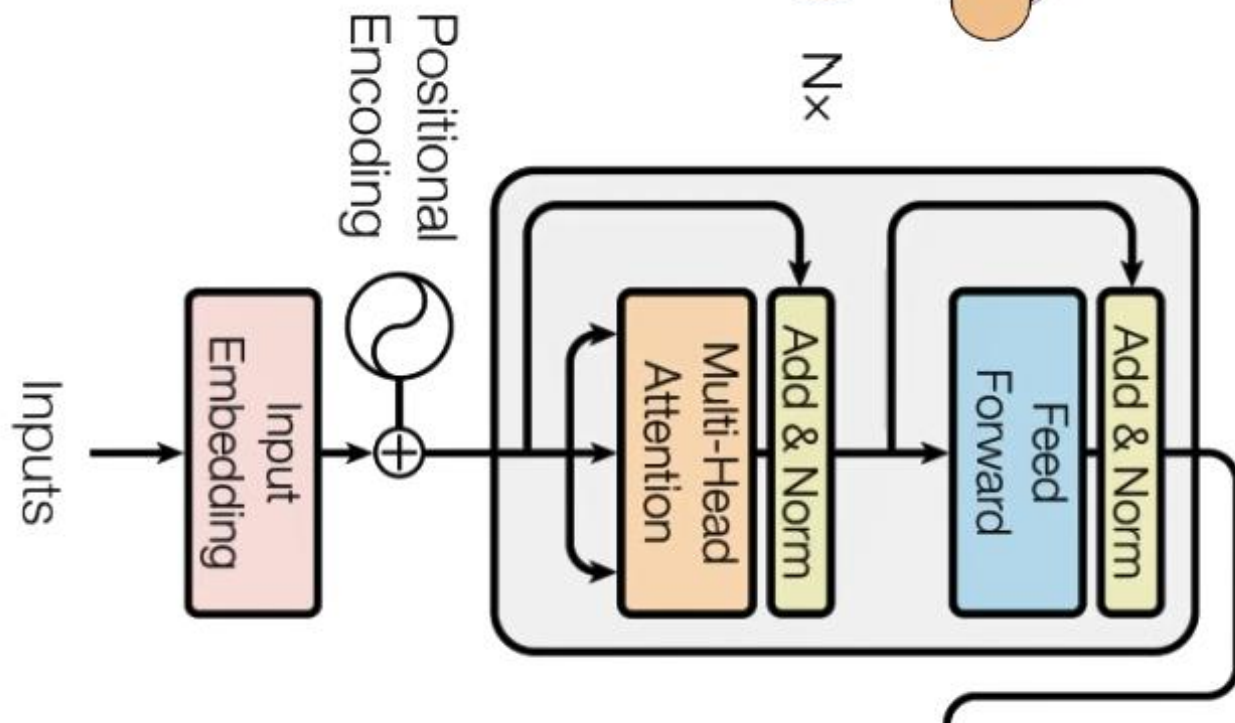
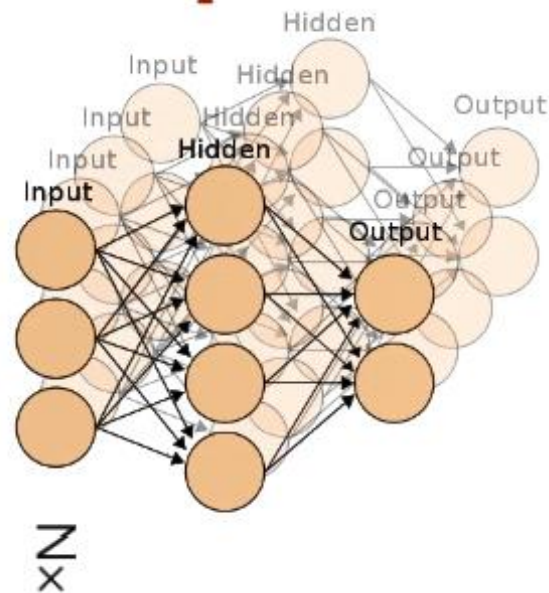
$[0.71 \quad 0.04 \quad 0.07 \quad 0.18]^T$
$[0.01 \quad 0.84 \quad 0.02 \quad 0.13]^T$
$[0.09 \quad 0.05 \quad 0.62 \quad 0.24]^T$
$[0.03 \quad 0.03 \quad 0.03 \quad 0.91]^T$

Transformer Components



The big red dog

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

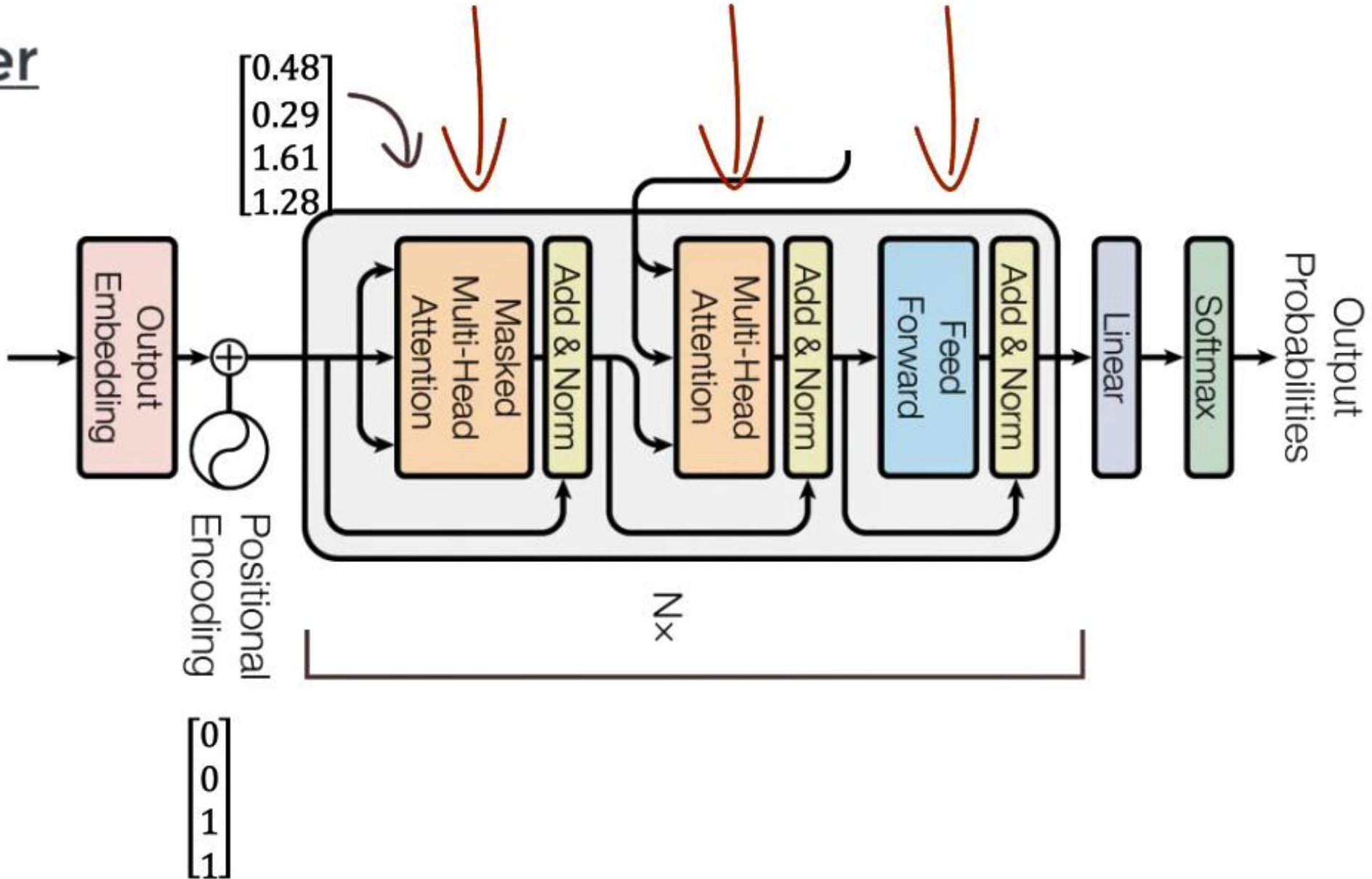


Transformer Components

Decoder

Le

Outputs
(shifted right)



Transformer Components

Decoder

Self
Attention

Le → Le gros chien rouge
gros → Le gros chien rouge
chien → Le gros chien rouge
rouge → Le gros chien rouge

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

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

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

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

Transformer Components

Decoder

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

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

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

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

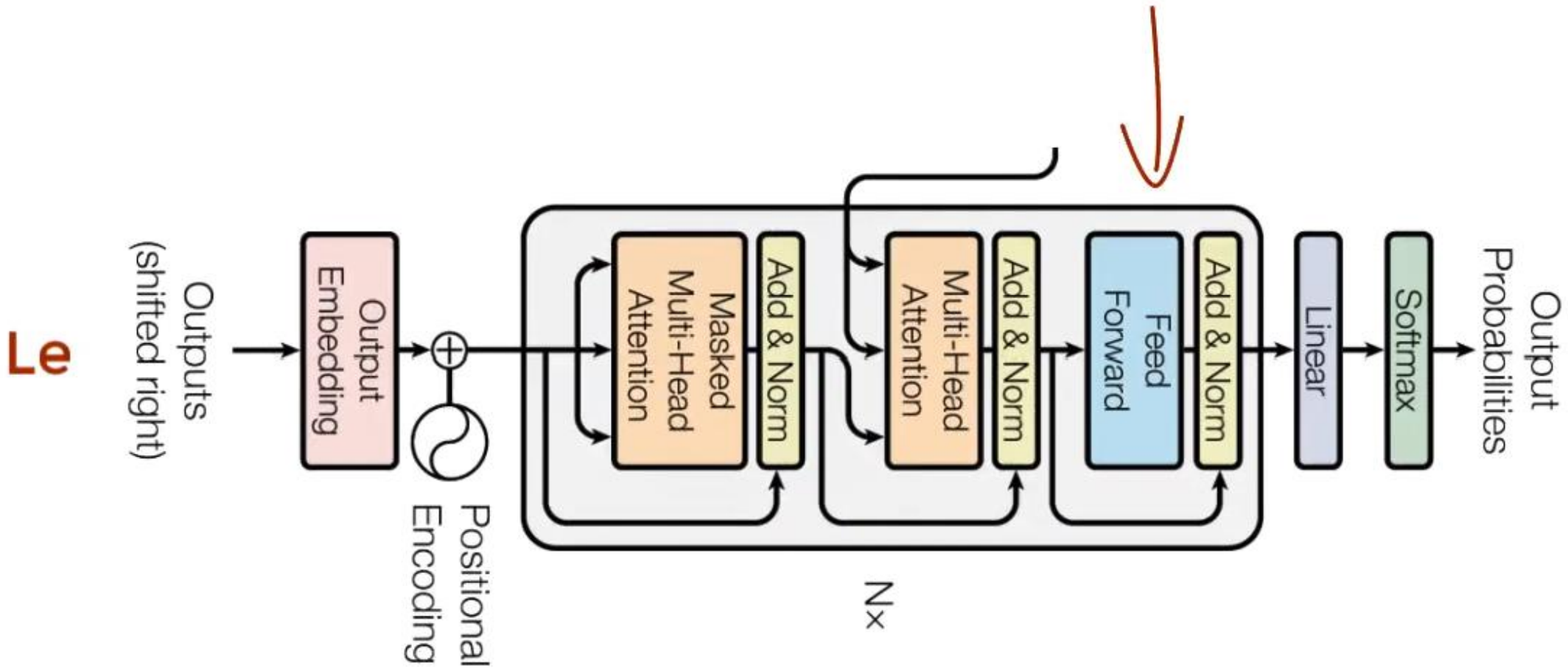
 $\begin{bmatrix} 0.03 \\ 0.03 \\ 0.03 \\ 0.91 \end{bmatrix}$

dog



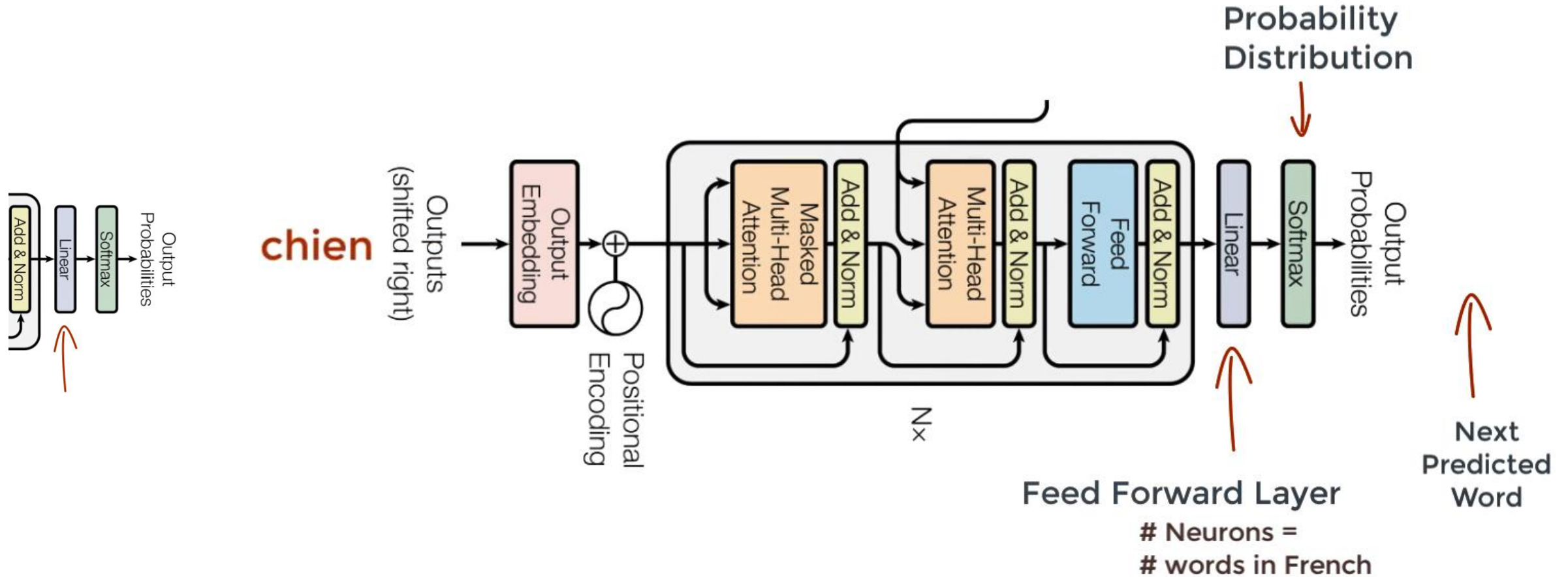
Encoder-
Decoder
Attention

Transformer Components



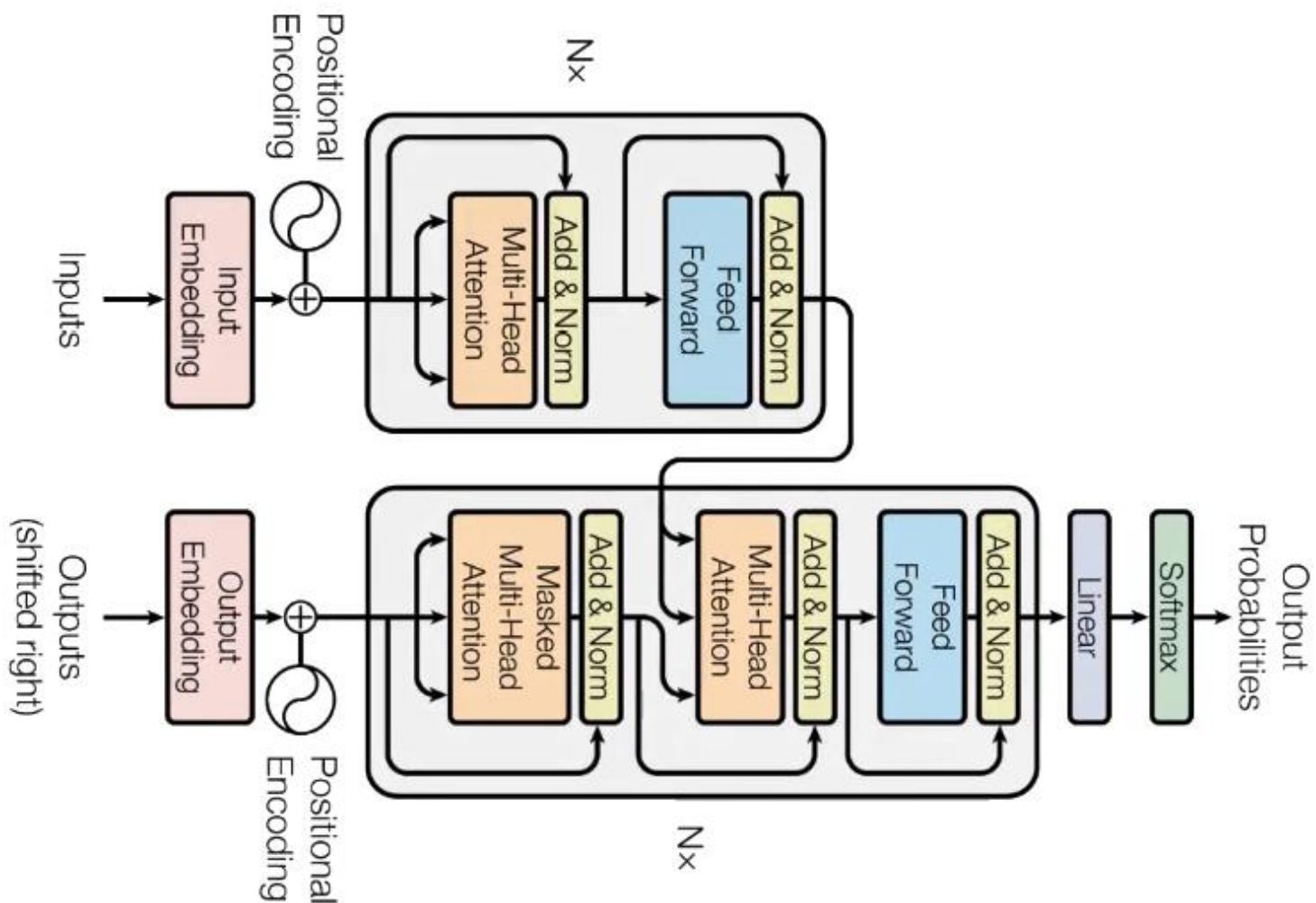
WE PASS EACH ATTENTION VECTOR TO THE FEED FORWARD UNIT

Transformer Components



Transformer Components

Next Pass



Multi-head Attention

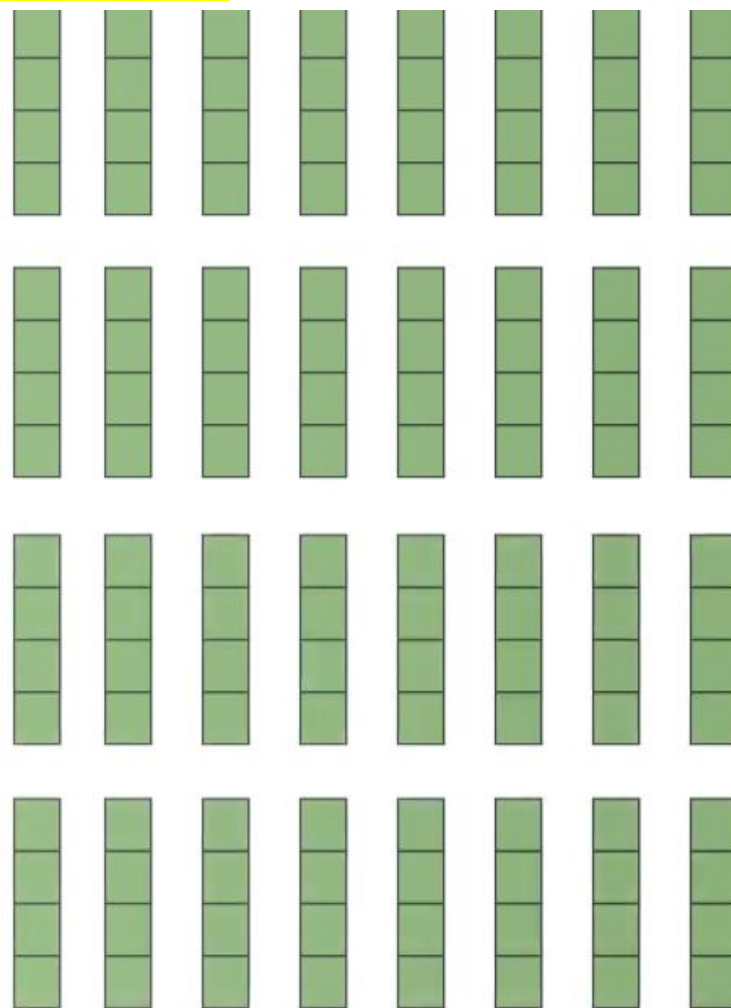
EIGHT ATTENTION VECTOR PER WORD!!

Attention

We want interactions
like this

Focus

The → The big red dog
big → The big red dog
red → The big red dog
dog → The big red dog



Averaged
vectors

The

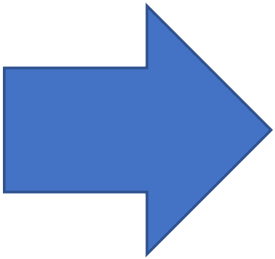
big

red

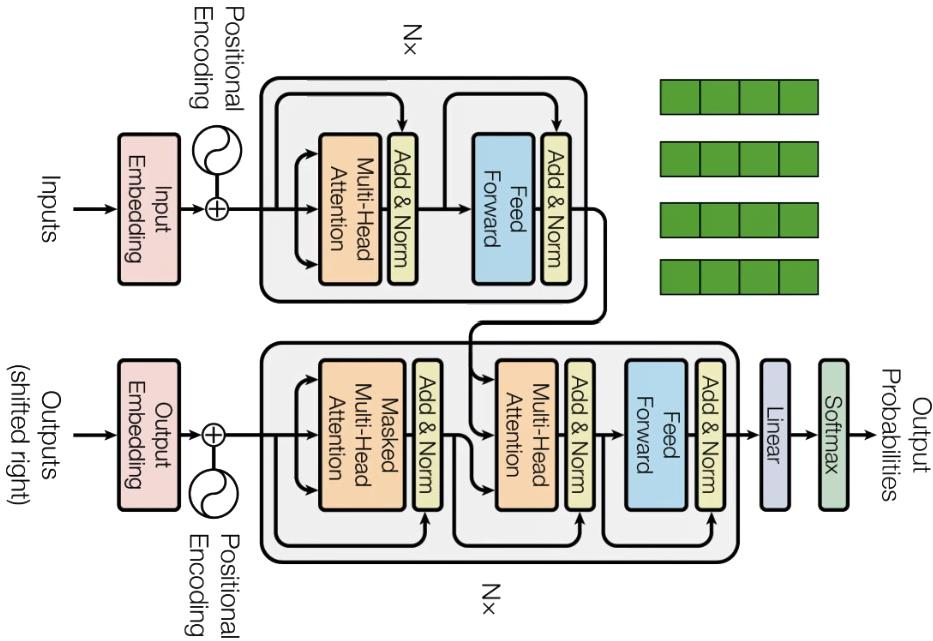
dog



The
big
red
dog



The
big
red
dog



Multi-headed Attention

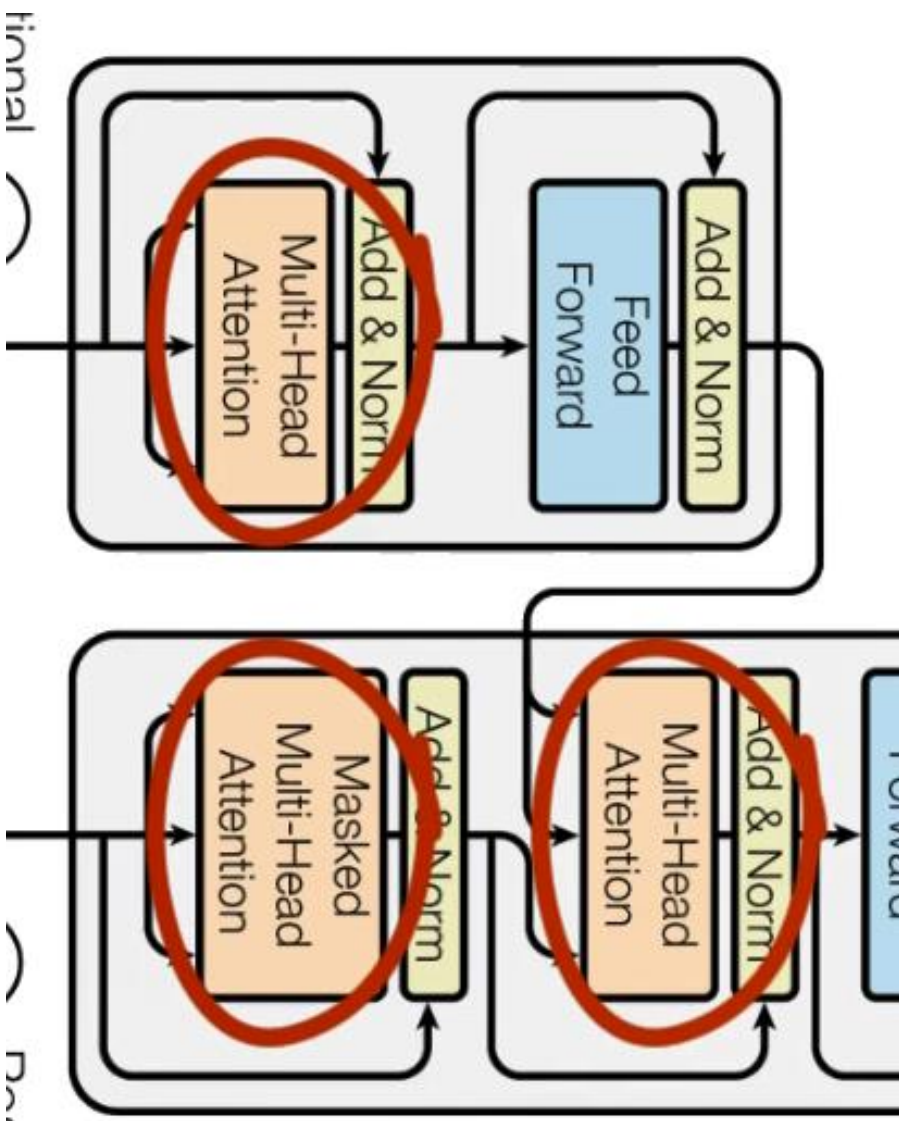
Encoder

The → The big red dog
big → The big red dog
red → The big red dog
dog → The big red dog

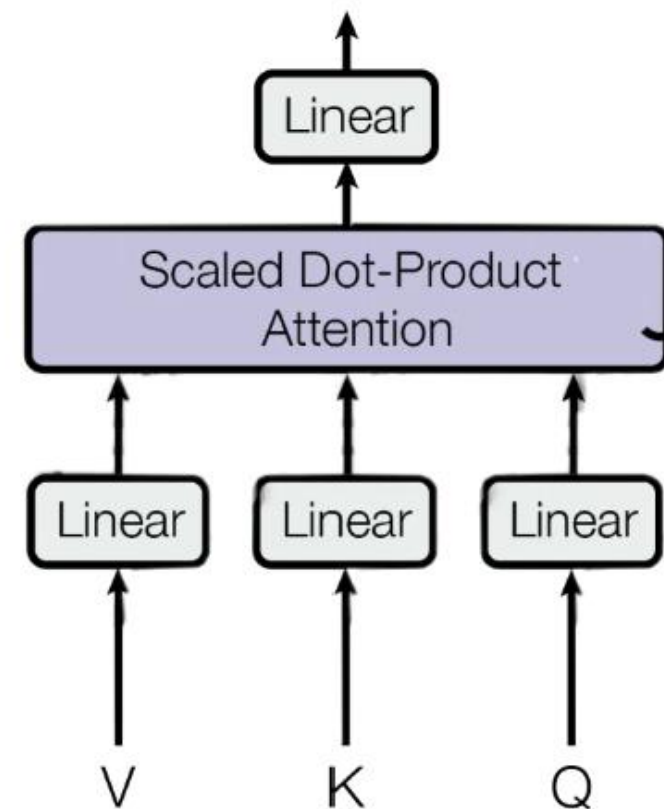
Decoder

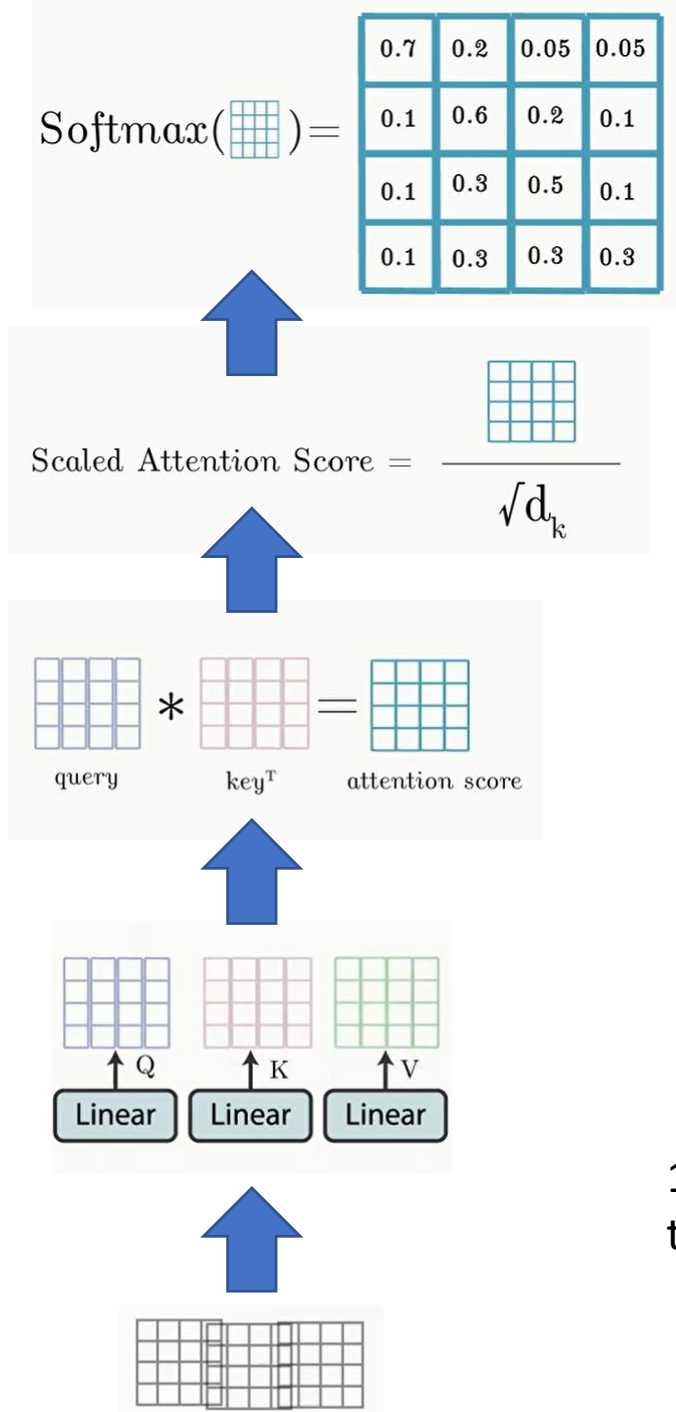
Le → Le gros chien rouge
gros → Le gros chien rouge
chien → Le gros chien rouge
rouge → Le gros chien rouge

Masked Input



**HOW MULTI-HEAD
ATTENTION LOOKS
LIKE ?**





5) These scores are passed into a soft max layer which returns probability values between 0 and 1. Finally this SOFTMAX attention scores are used to multiply the value vector to produce a weighted output vector

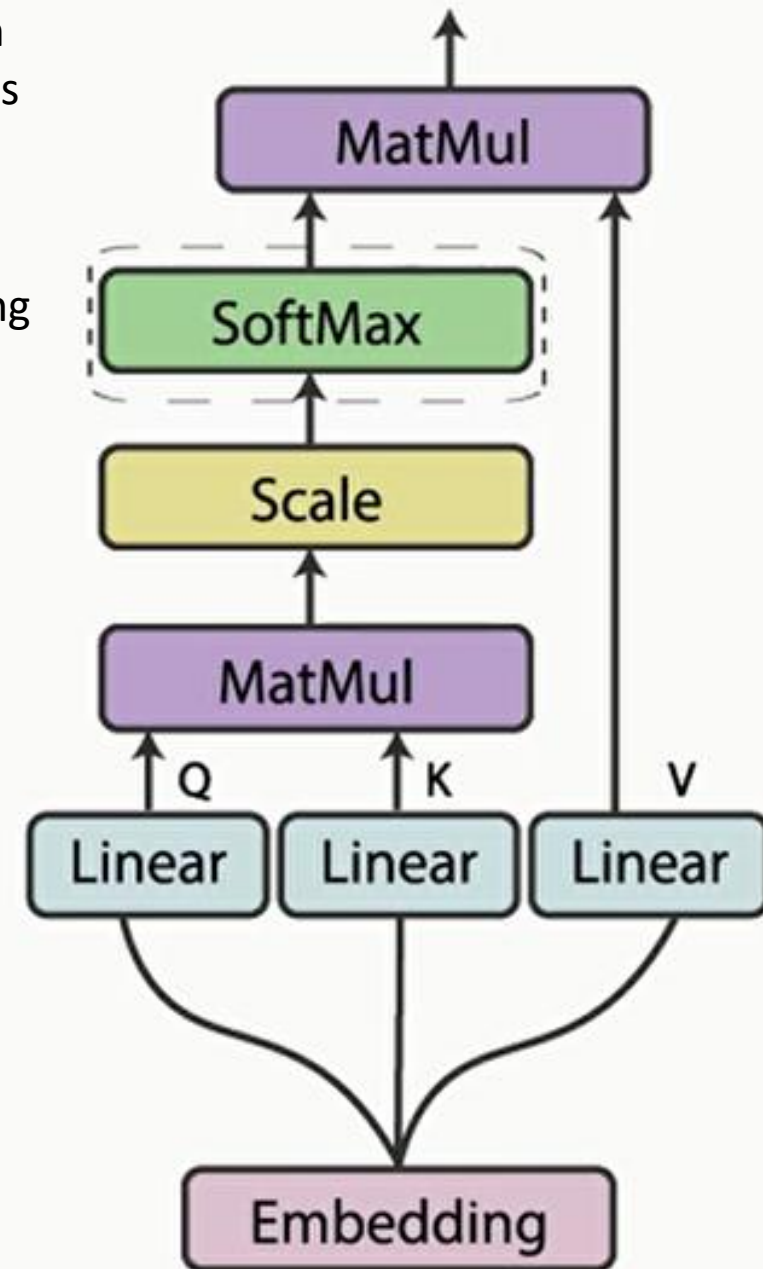
4) The scoring matrix then gets scaled down by dividing by the square root of the dimension of the key vector and this is to allow for more stable gradients as multiplying these values several times can have exploding effects

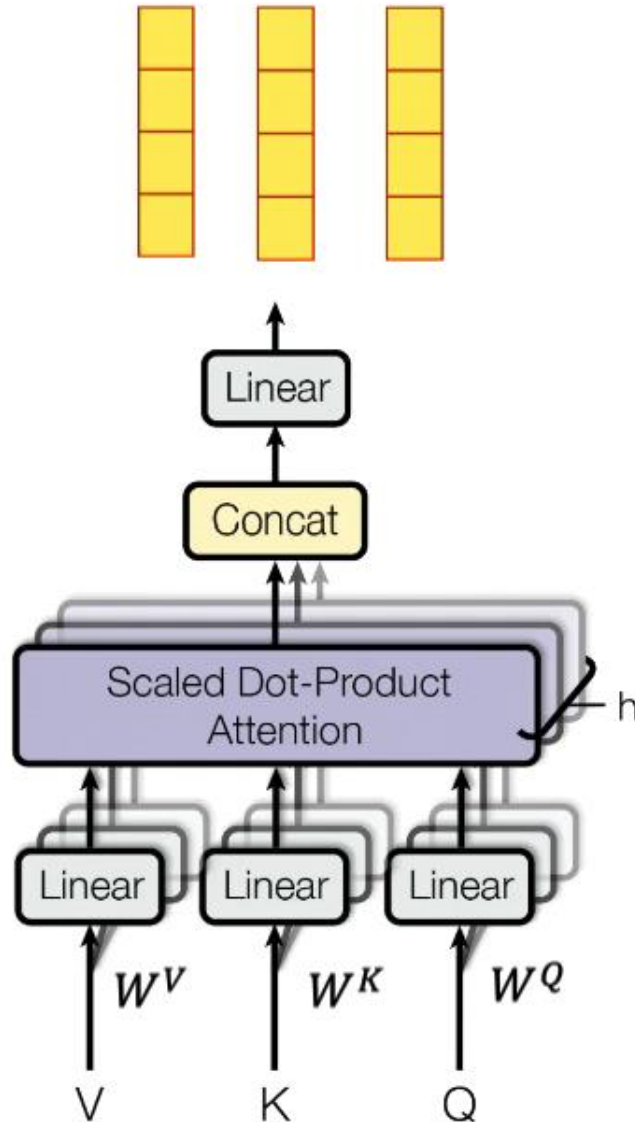
3) Dot multiplication is applied to the query and the transpose of the key to produce a scoring matrix and this scoring matrix is what determines the relevance of a word to other words in the matrix the higher the score the more relevant the word is

2) Create the ENCODER attention query Q key K and value V vectors

1) The positional input embeddings are feed into three distinct linear layers

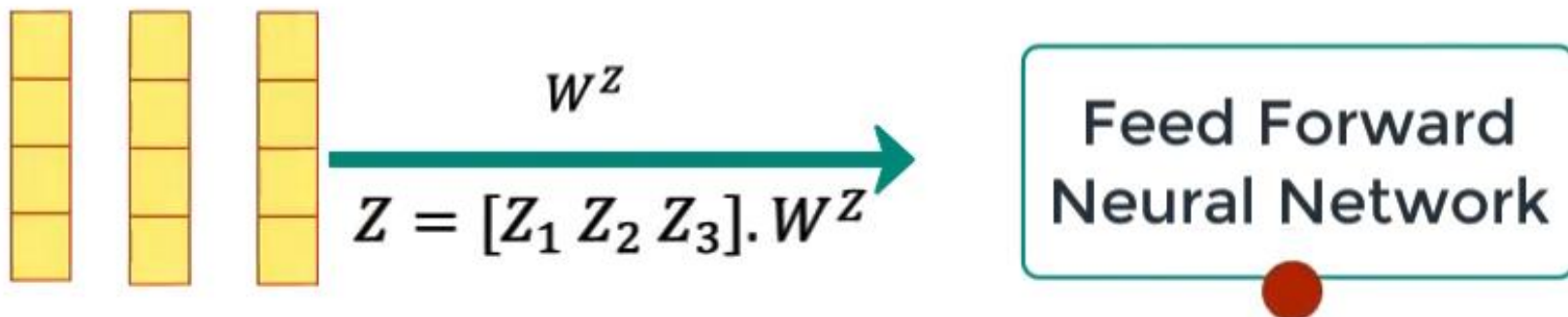
$$Attention(Q, K, V) = softmax_k \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$





MULTIPLE ATTENTION Z for every word!!

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

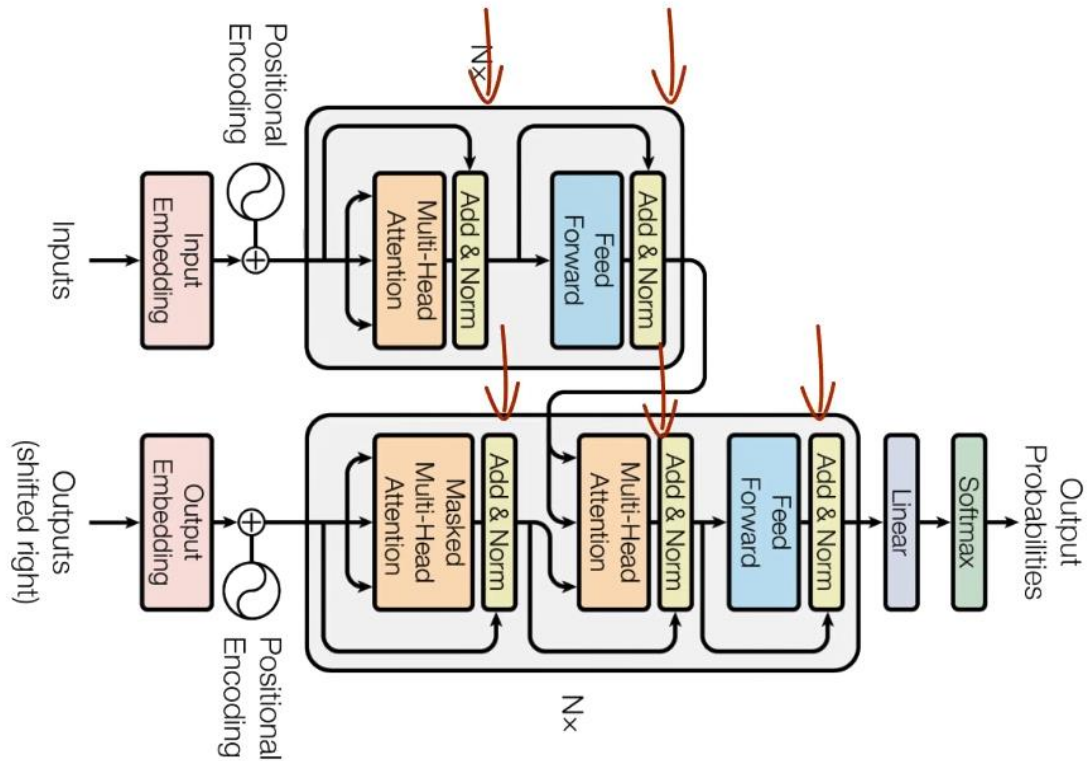


**Matrix $Wz \rightarrow$ OUTPUT “Z”
IS STILL ATTENTION
VECTOR PER WORD**

$$Z = \text{softmax} \left(\frac{Z \cdot V}{\sqrt{d_v}} \right) \cdot V$$

**Nice
combo**

AFTER MULTI-HEAD ATTENTION – WE ADD THE NORMALIZATION STEP



Batch Normalization

batch

			Same for all training examples	
			mean	std
1	3	6	3	3
2	2	2	2	0
0	1	5	3	3
4	6	1	4	3
5	2	3	3	2
1	0	1	1	1

