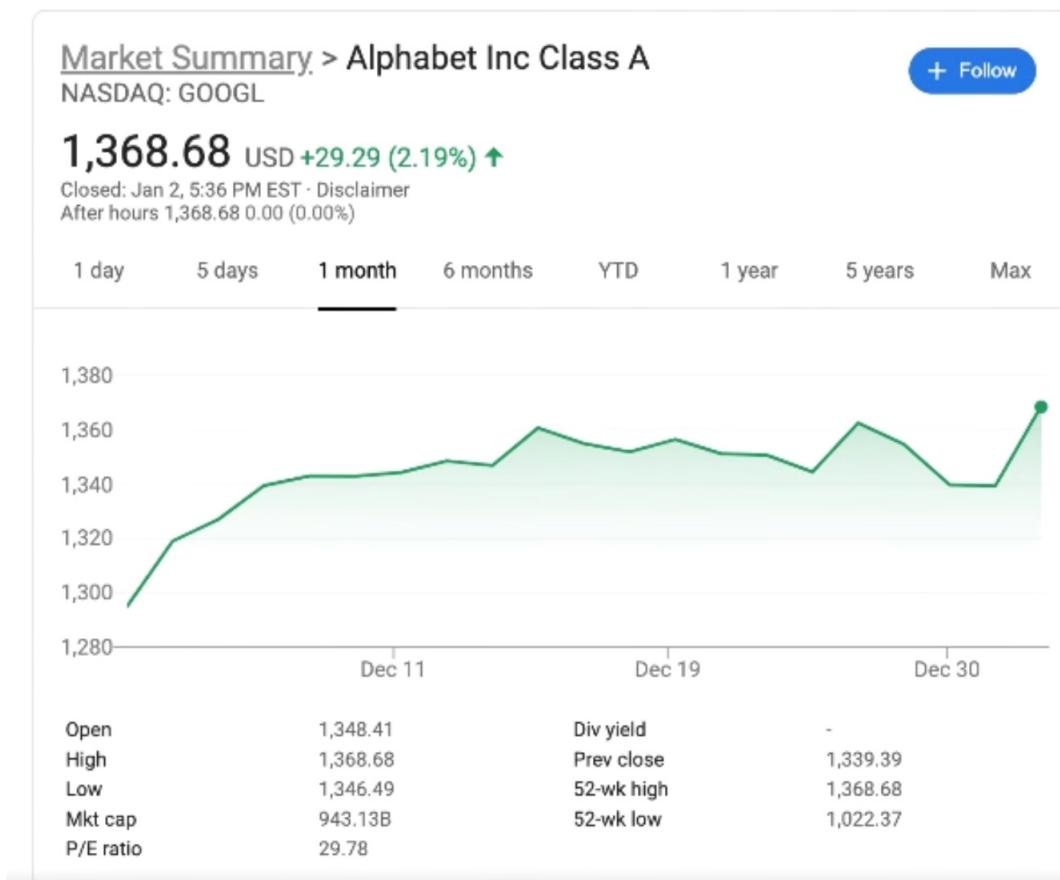


Recurrent Neural Networks

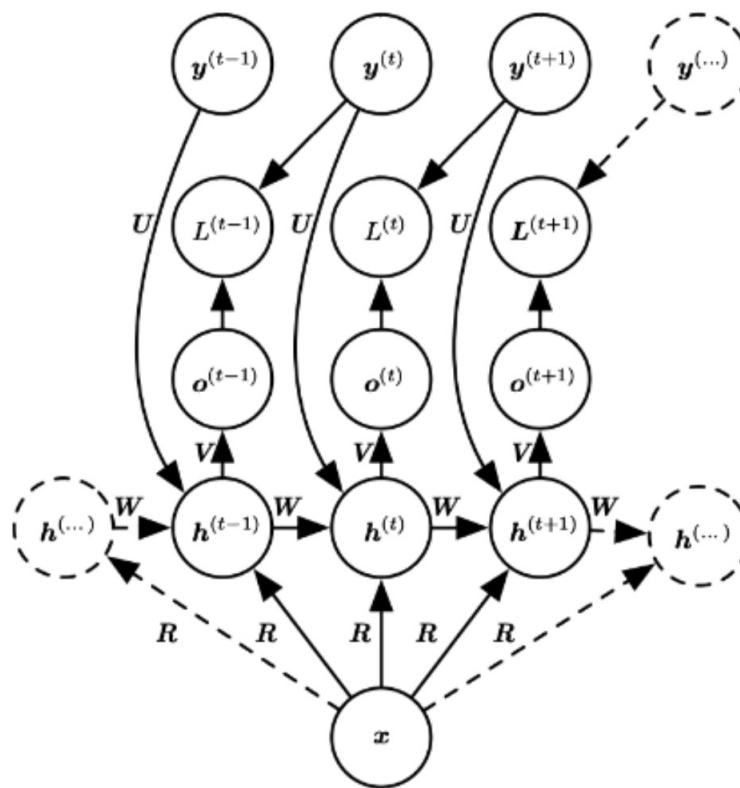
The unicorn is scotland's national animal

Sequences



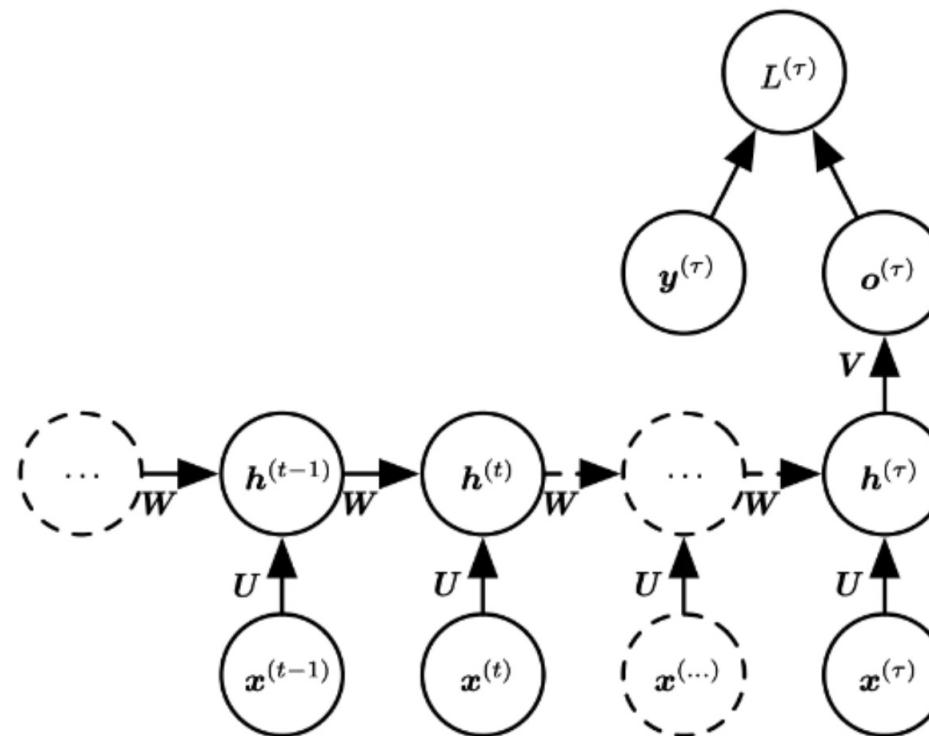
Recurrent Neural Networks

1. Vector-Sequence Models



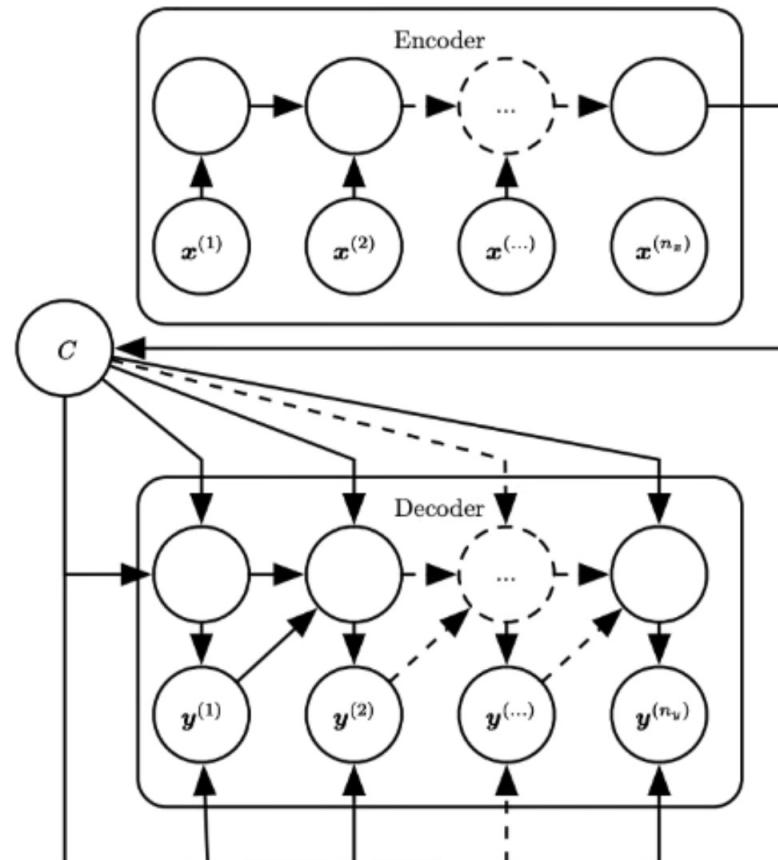
Recurrent Neural Networks

2. Sequence-Vector Models



Recurrent Neural Networks

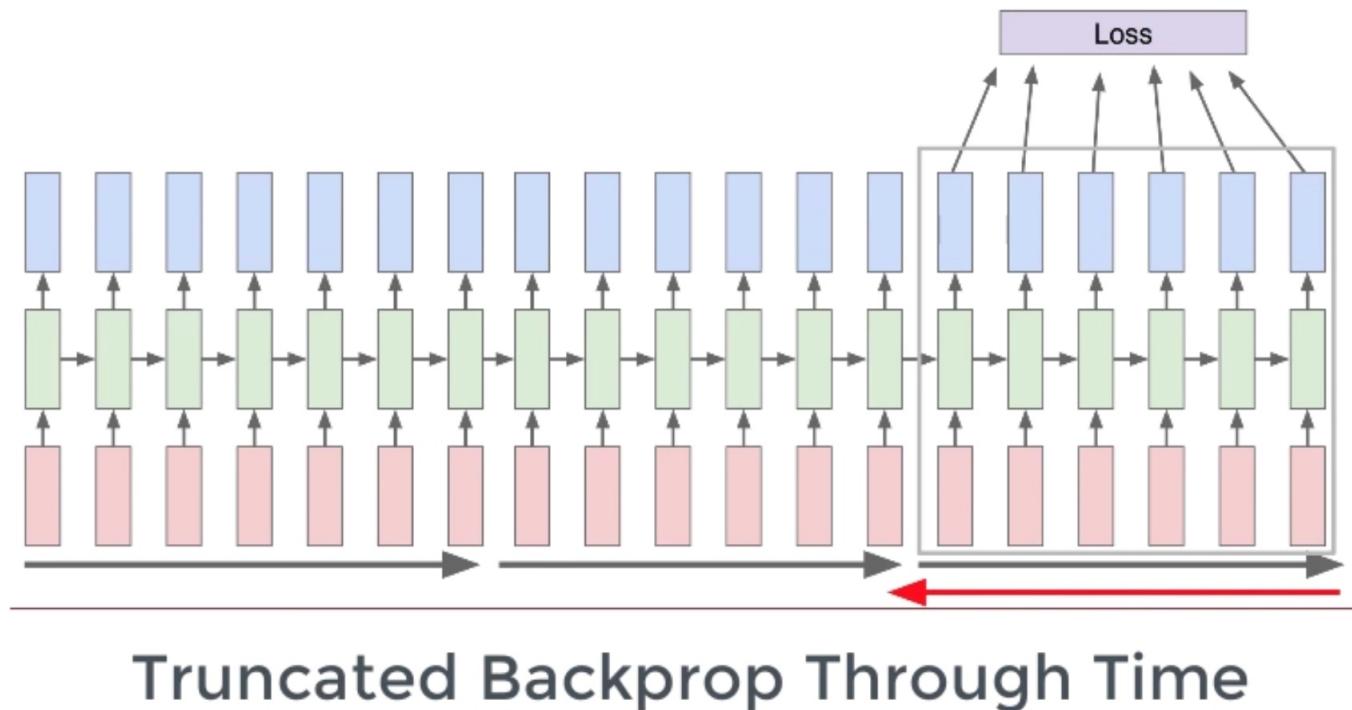
3. Sequence-Sequence Models



Recurrent Neural Networks

Disadvantages

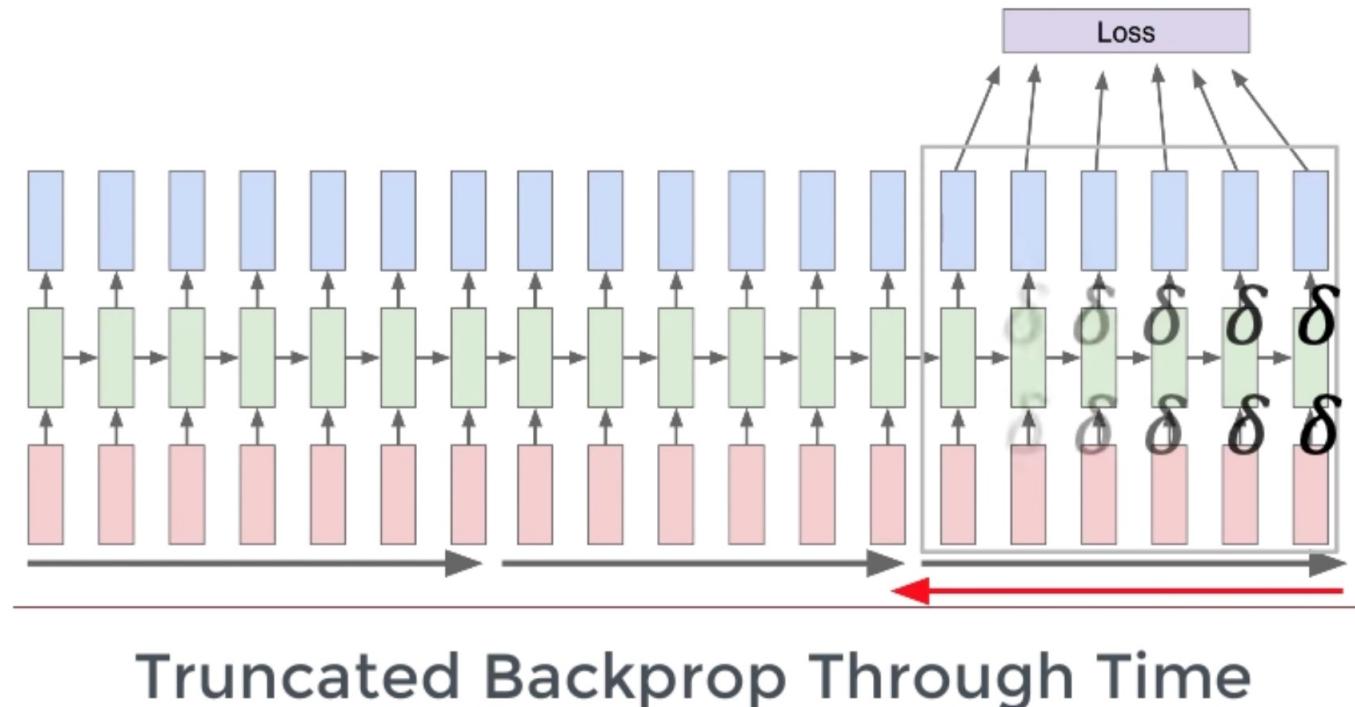
1. Slow to train.



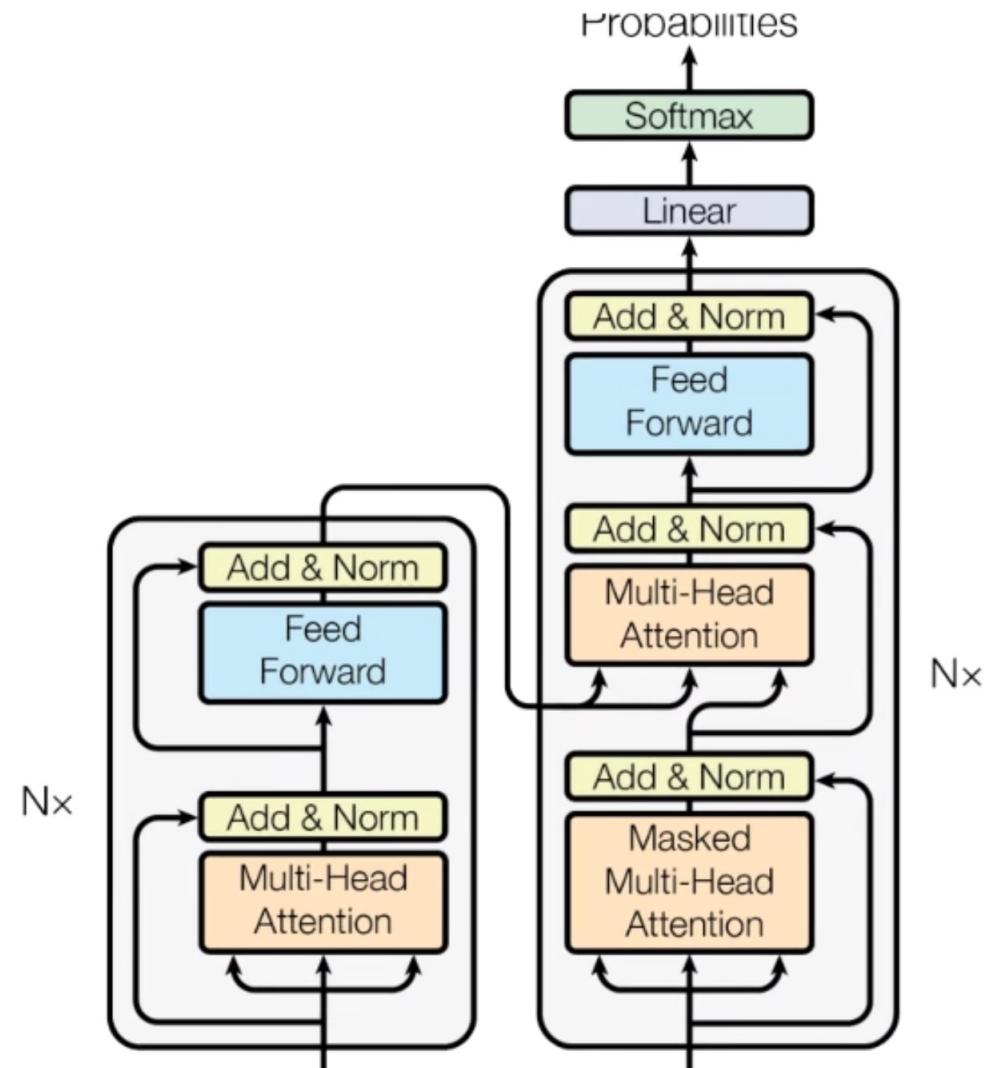
Recurrent Neural Networks

Disadvantages

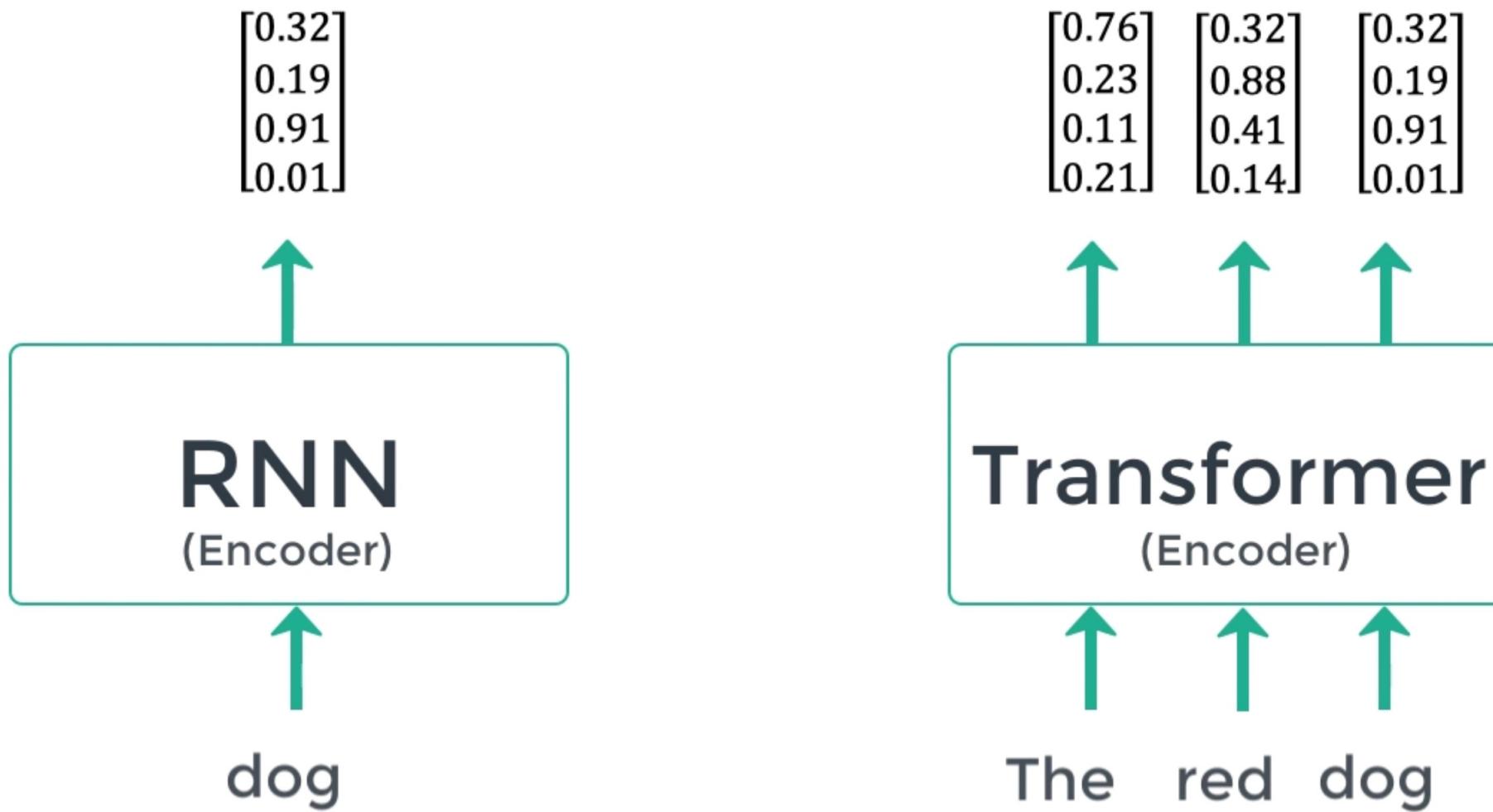
1. Slow to train.
2. Long sequences lead to vanishing/exploding gradients



Transformers



English-French Translation



Transformer Components

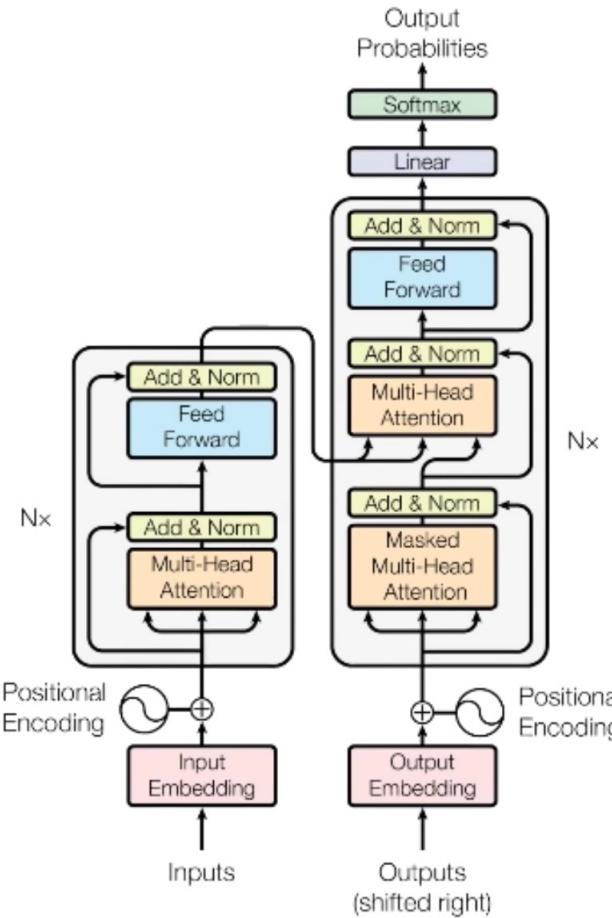
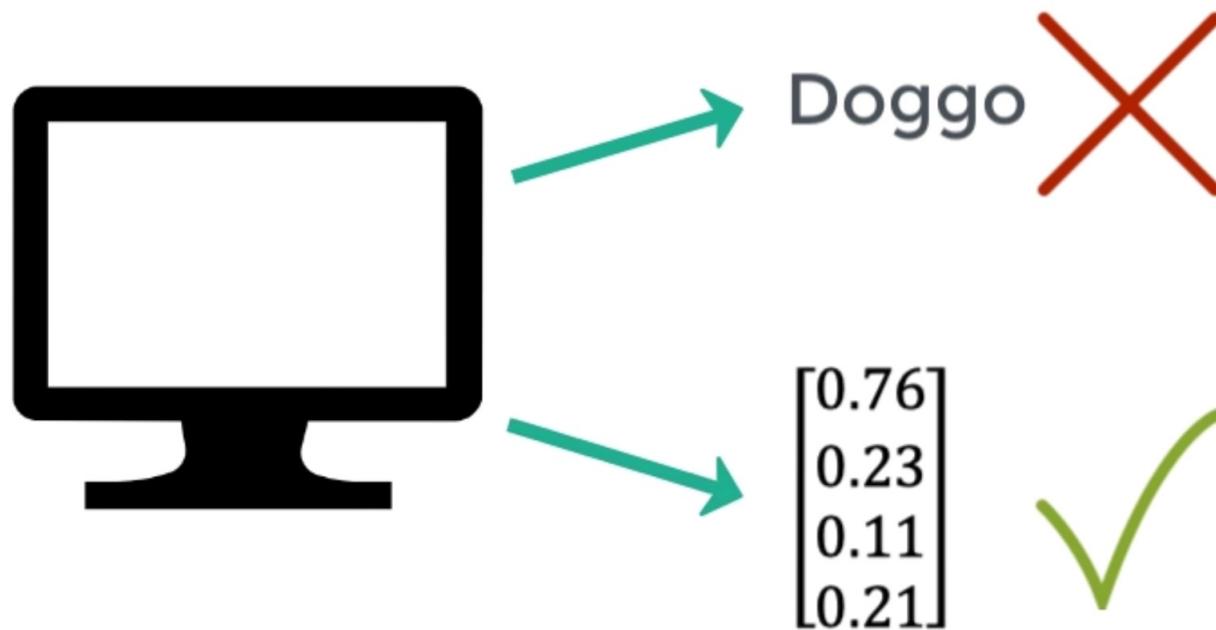


Figure 1: The Transformer - model architecture.

Transformer Components

Input Embedding



Transformer Components

Input Embedding



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](#)



SUBSCRIBE

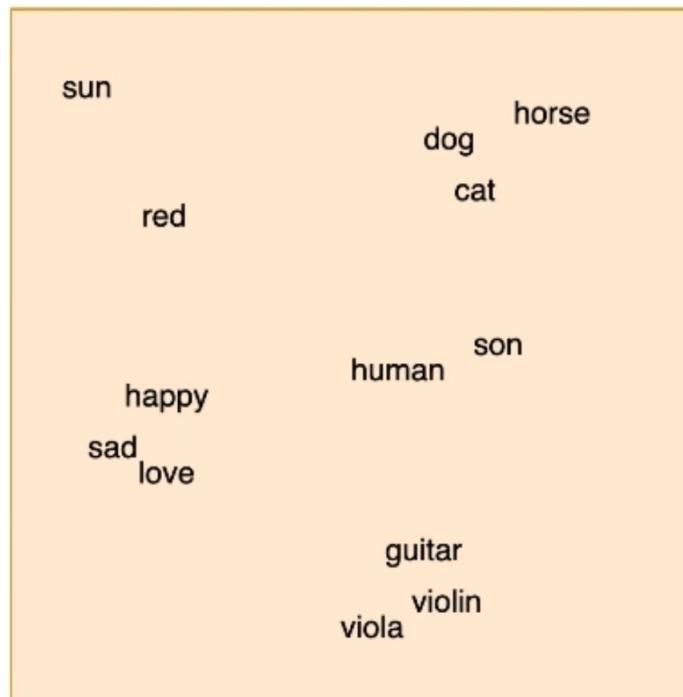
Transformer Components

Input Embedding

dog 

AJ's **dog** is a cutie

AJ looks like a **dog**



[0.37
0.99
0.01
0.08]

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

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

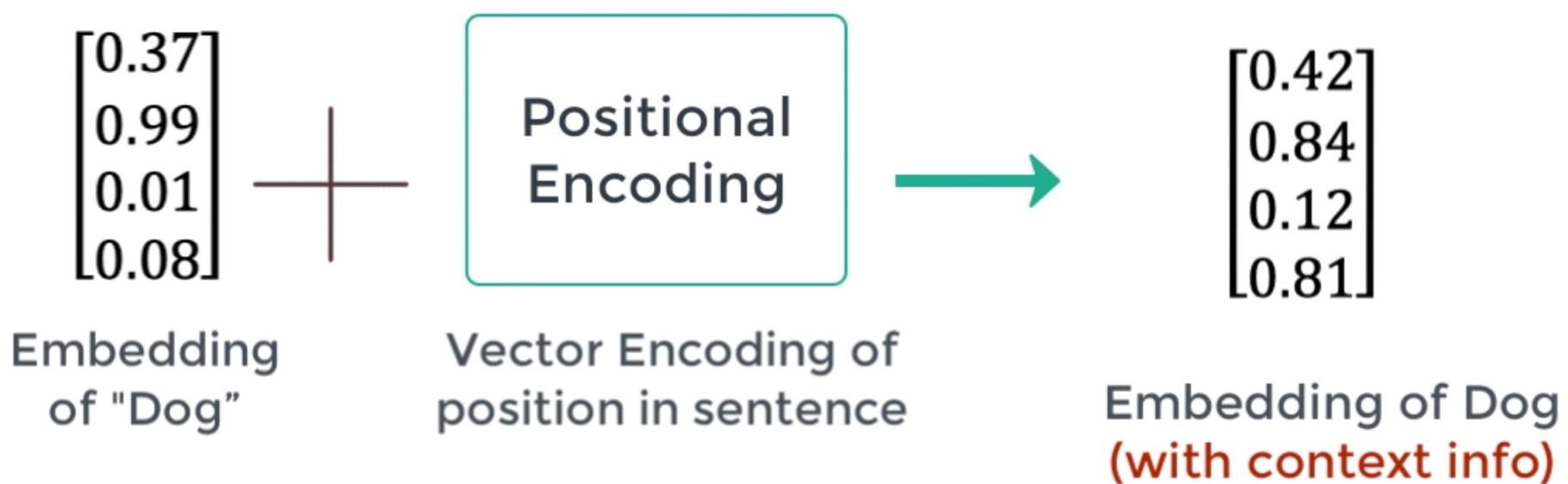
$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

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

Transformer Components

Positional Encoder

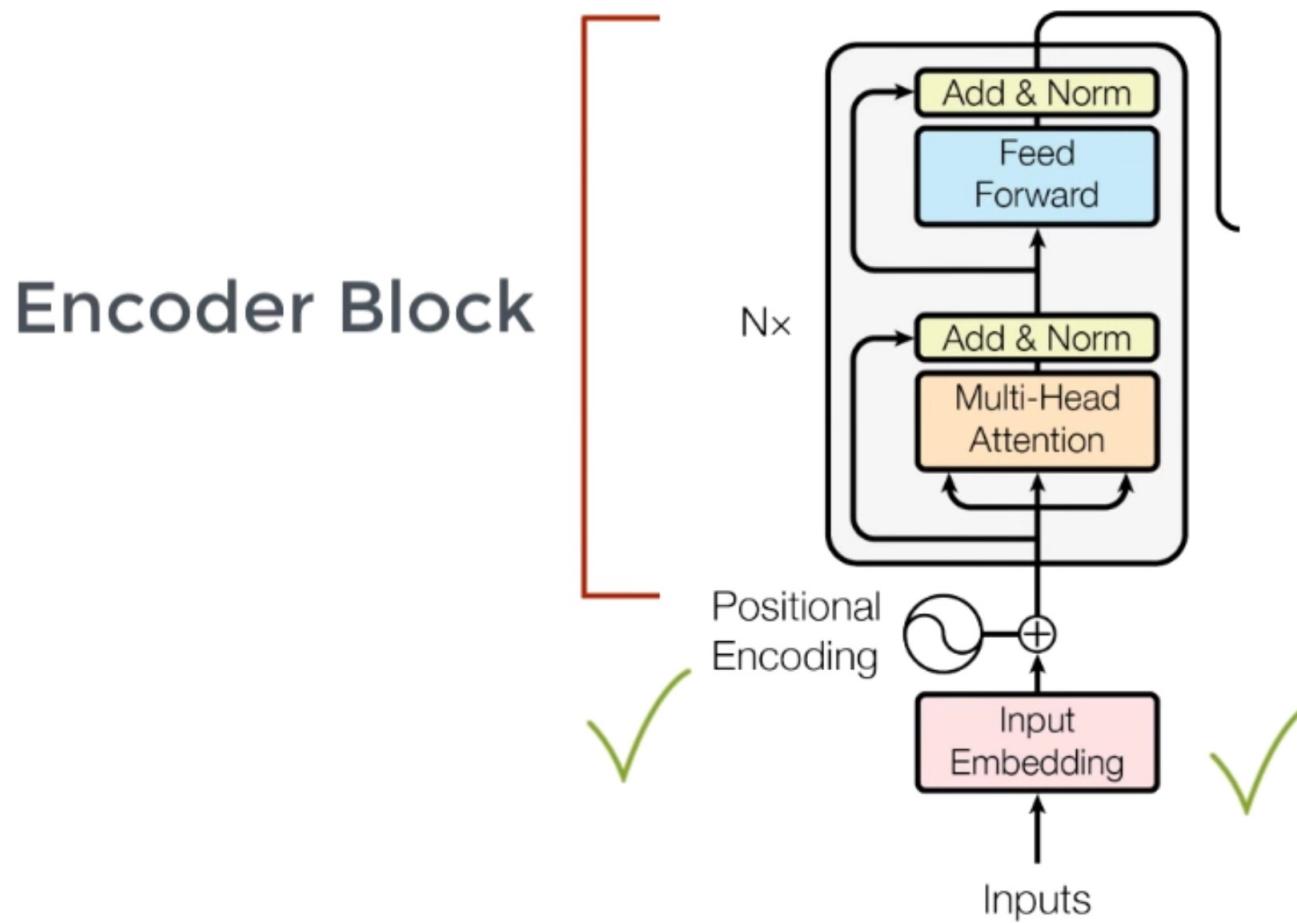
:vector that gives context based on position of word in sentence



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

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

Transformer Components



Transformer Components

Attention : What part of the input should we focus?

Transformer Components

Attention : What part of the input should we focus?

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



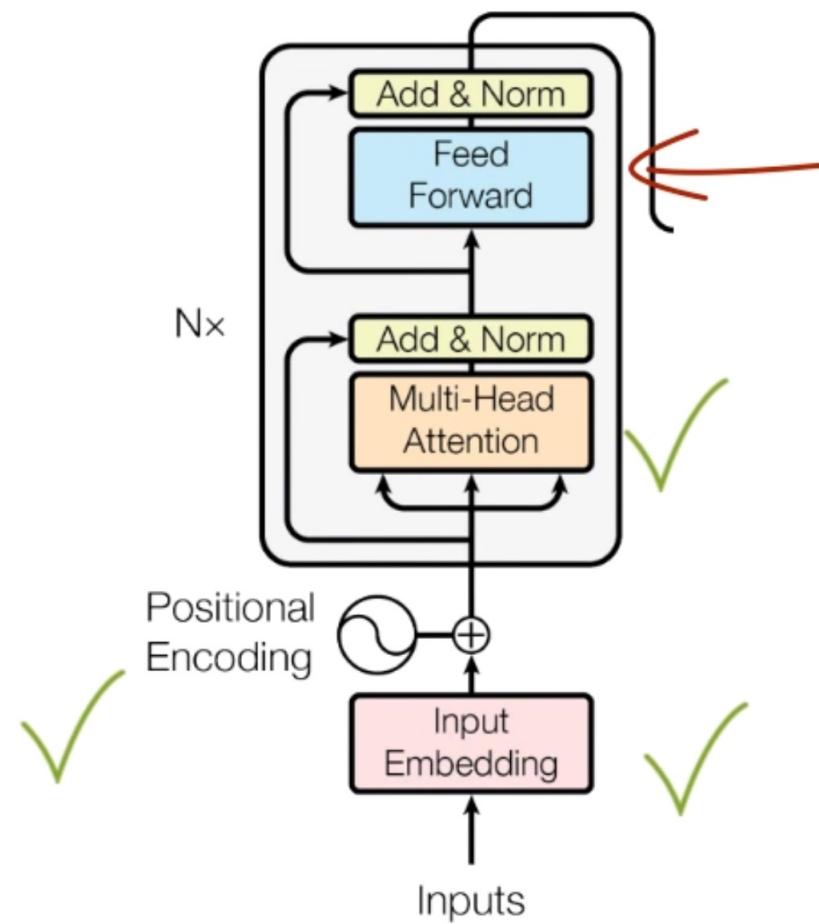
Focus

Transformer Components

Attention : What part of the input should we focus?

| | | Attention Vectors |
|-----|---------|------------------------------------|
| The | Focus → | [0.71 0.04 0.07 0.18] ^T |
| big | → | [0.01 0.84 0.02 0.13] ^T |
| red | → | [0.09 0.05 0.62 0.24] ^T |
| dog | → | [0.03 0.03 0.03 0.91] ^T |

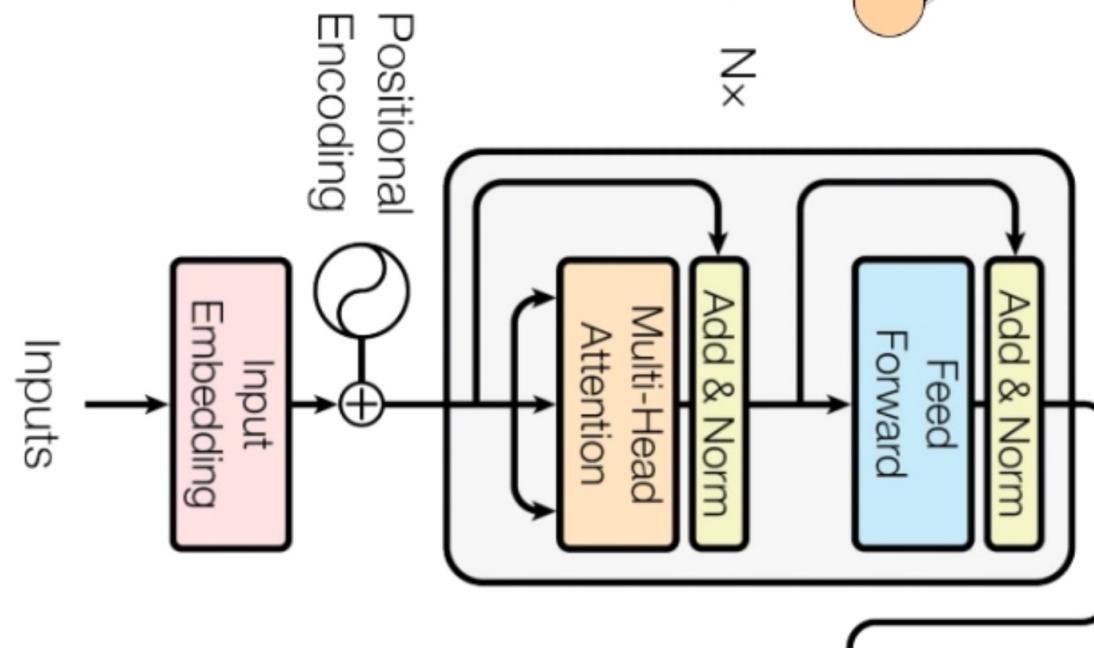
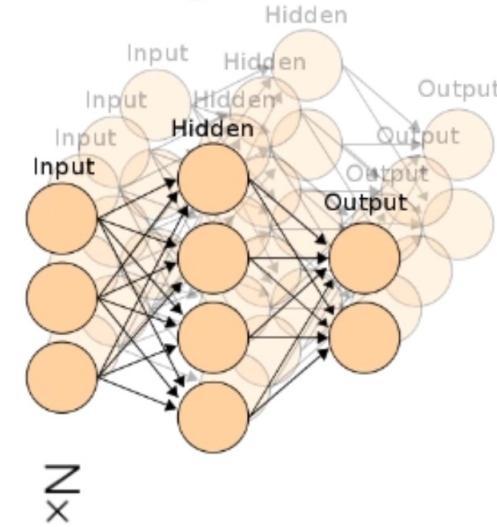
Transformer Components



Transformer Components

dog
red
big
The

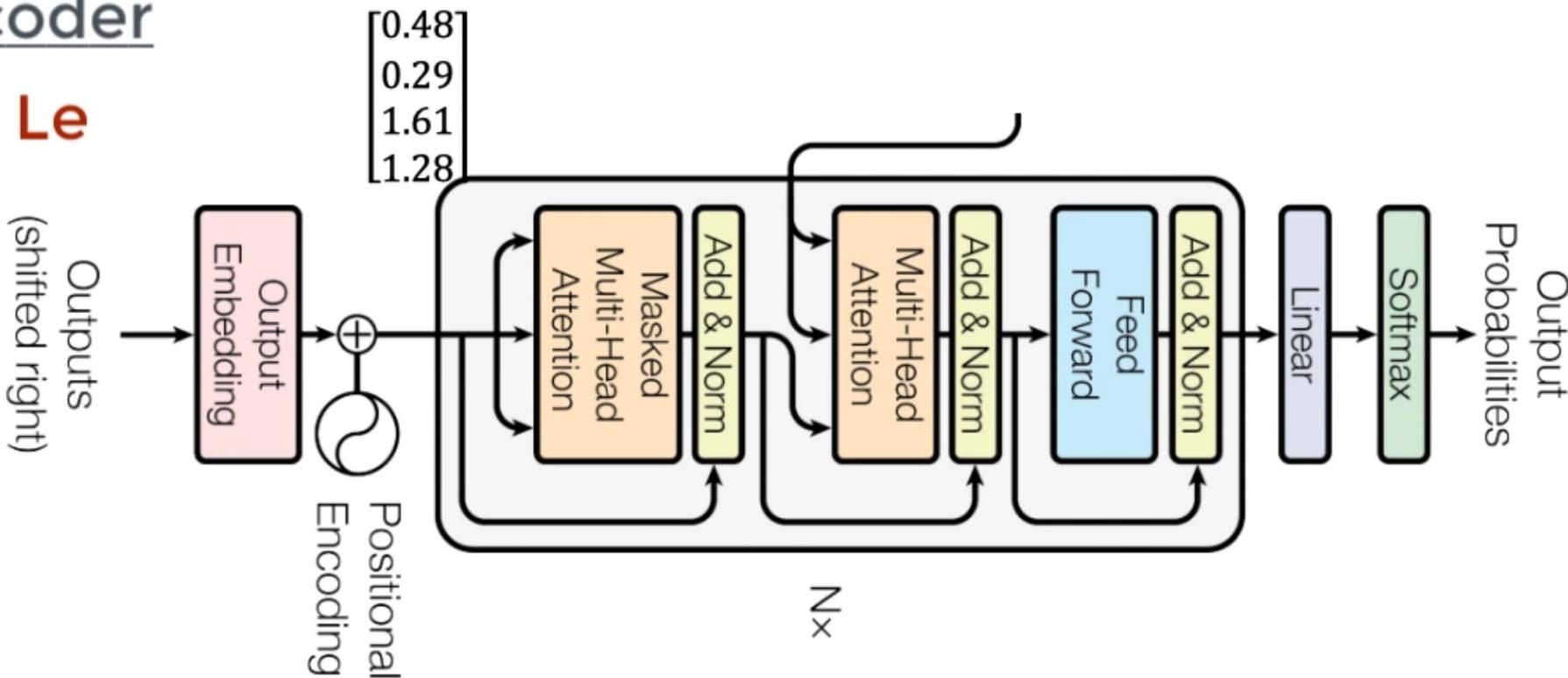
$$\begin{bmatrix} 0.03 \\ 0.09 \\ 0.03 \\ 0.03 \\ 0.71 \\ 0.04 \\ 0.07 \\ 0.18 \end{bmatrix}$$



Transformer Components

Decoder

Le



Outputs
(shifted right)

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

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

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

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

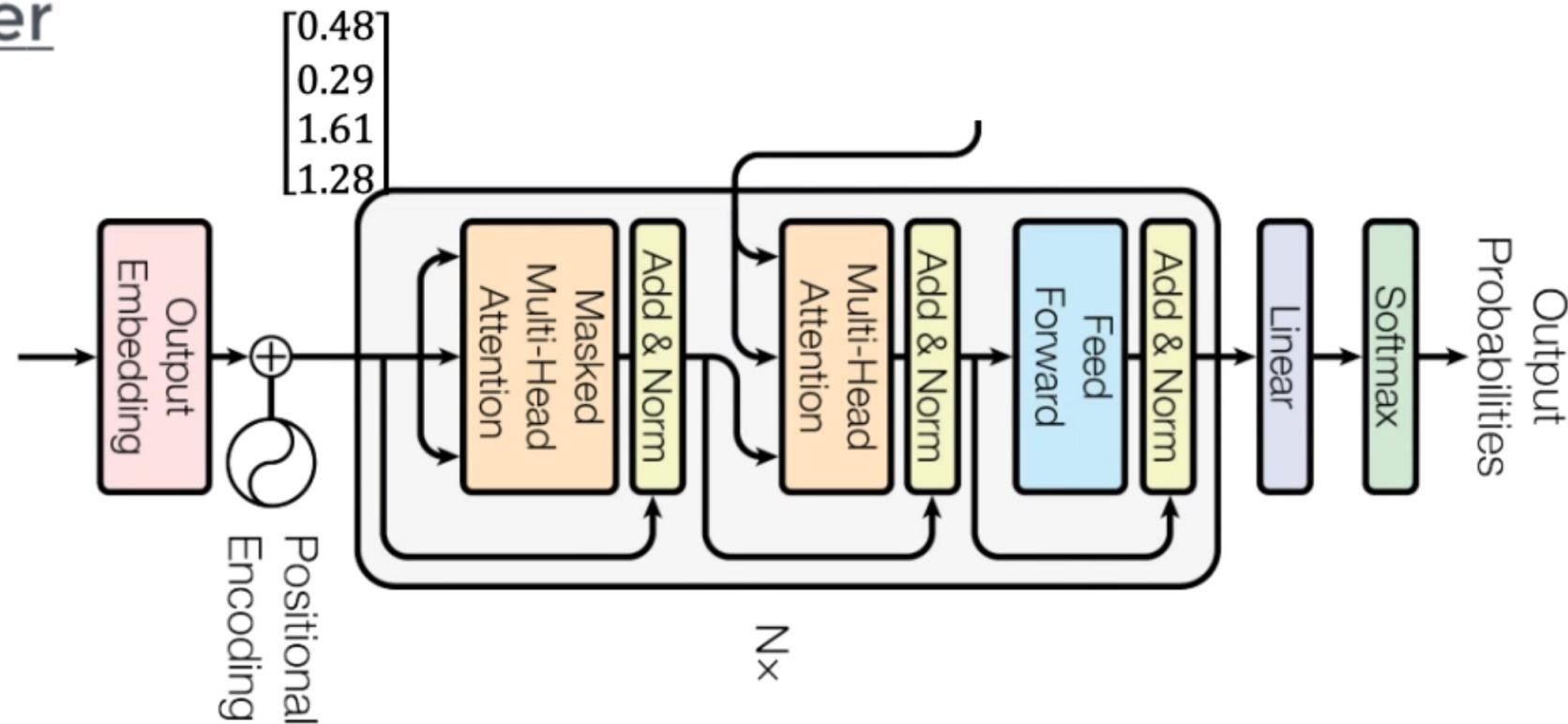


Transformer Components

Decoder

Le

Outputs
(shifted right)



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

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

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

Transformer Components

Decoder

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

Multi-headed Attention

Encoder

Decoder

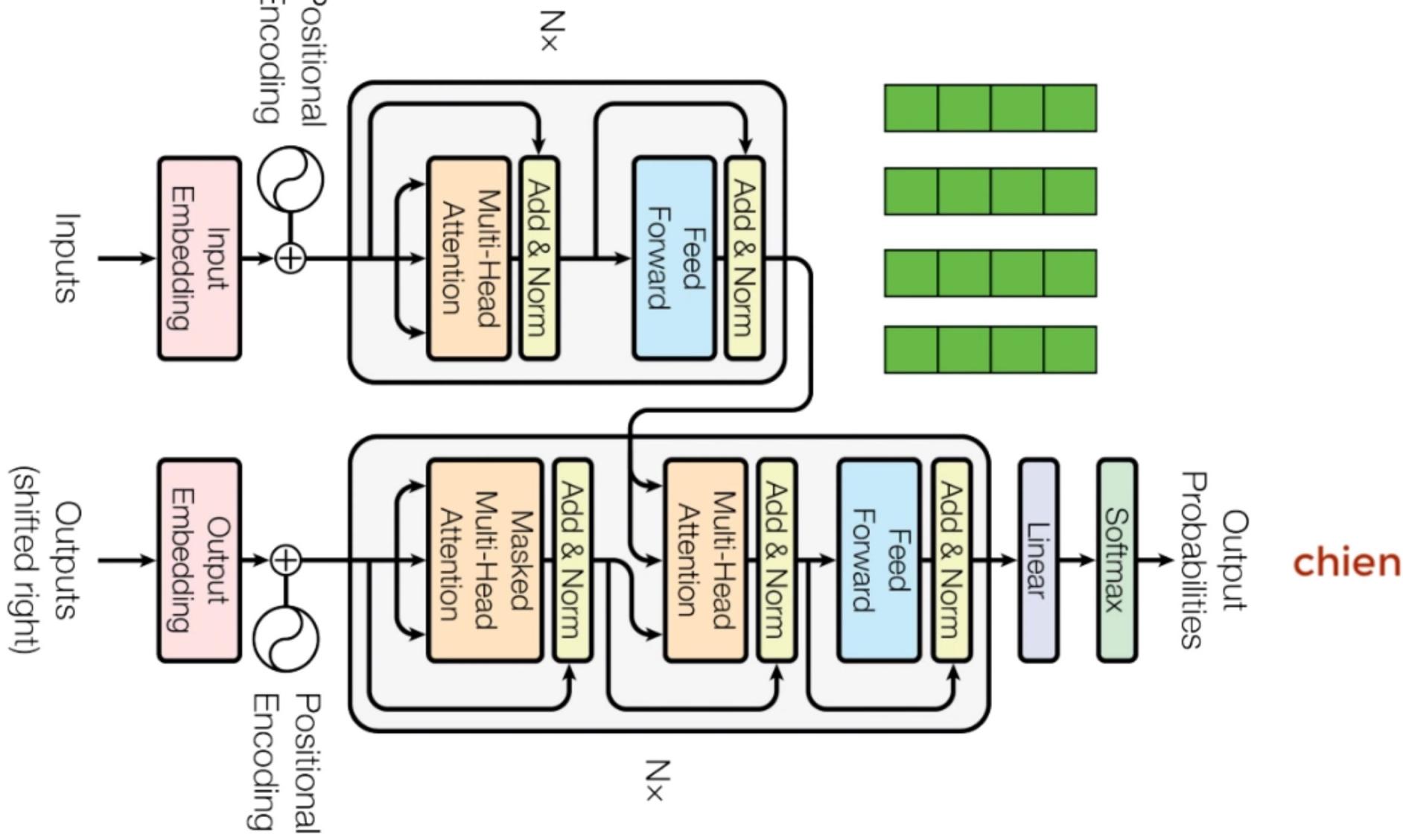
The → The big red dog
big → The big red dog
red → The big red dog
dog → The big red dog

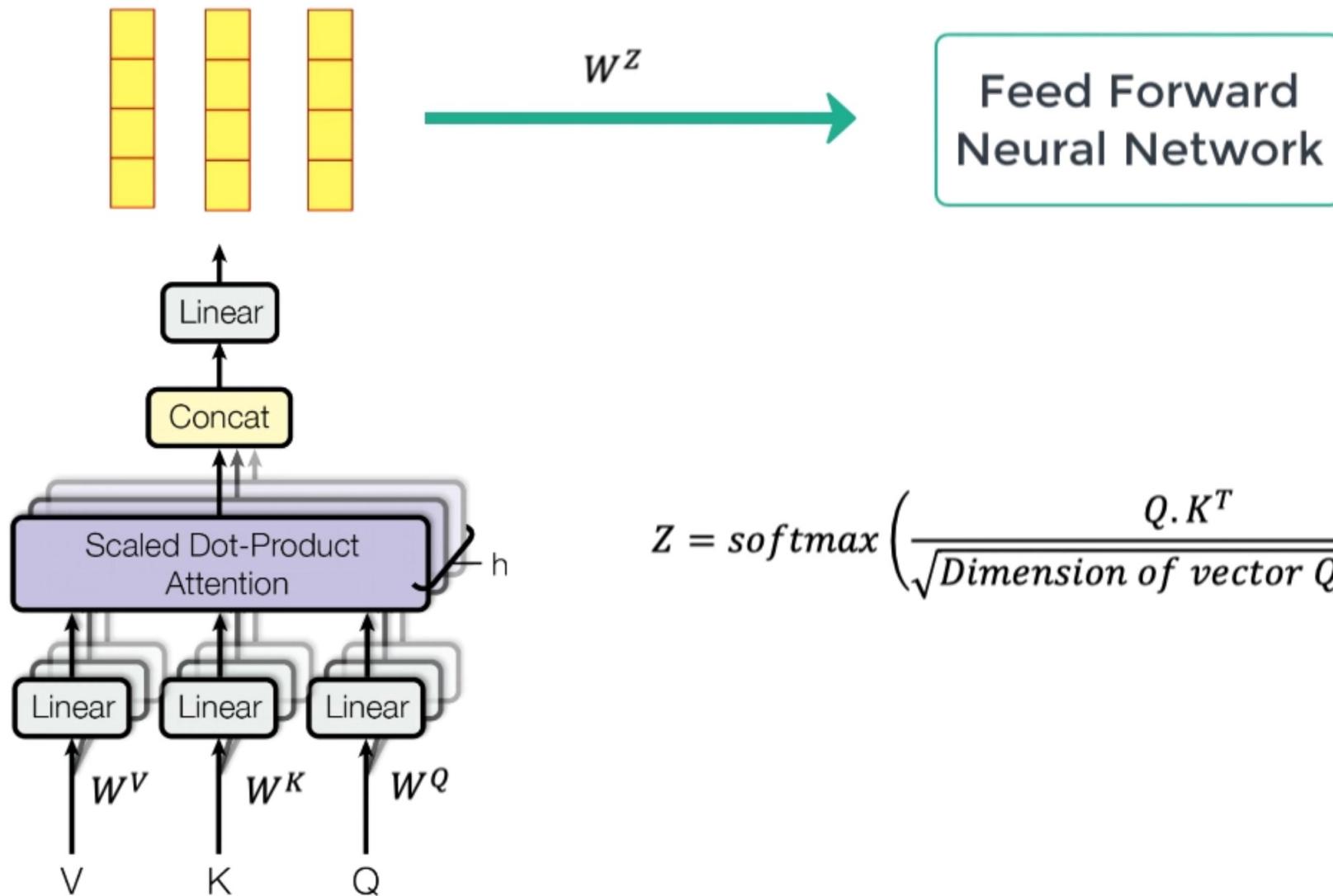
Le → Le gros chien rouge
gros → Le gros chien rouge
chien → Le gros chien rouge
rouge → Le gros chien rouge

Masked Input

The
big
red
dog

Le





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

Transformer Components

