



1

The image shows a slide titled 'Outline' in red. The outline contains the following points:

- Recurrent Neural Networks (RNNs)
  - Motivation, basic concepts, examples
- LSTM & GRU
- Transformer
  - Example applications

In the top right corner of the slide area, there is a small version of the Swinburne University of Technology logo, identical to the one on the title slide.

2

## Sequences & Sequence modelling

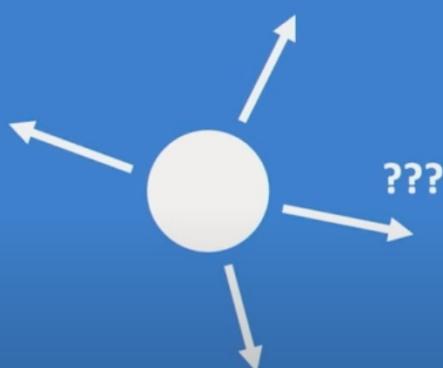
Given an image of a ball,  
can you predict where it will go next?



Source: MIT Introduction to Deep Learning 6.S191

3

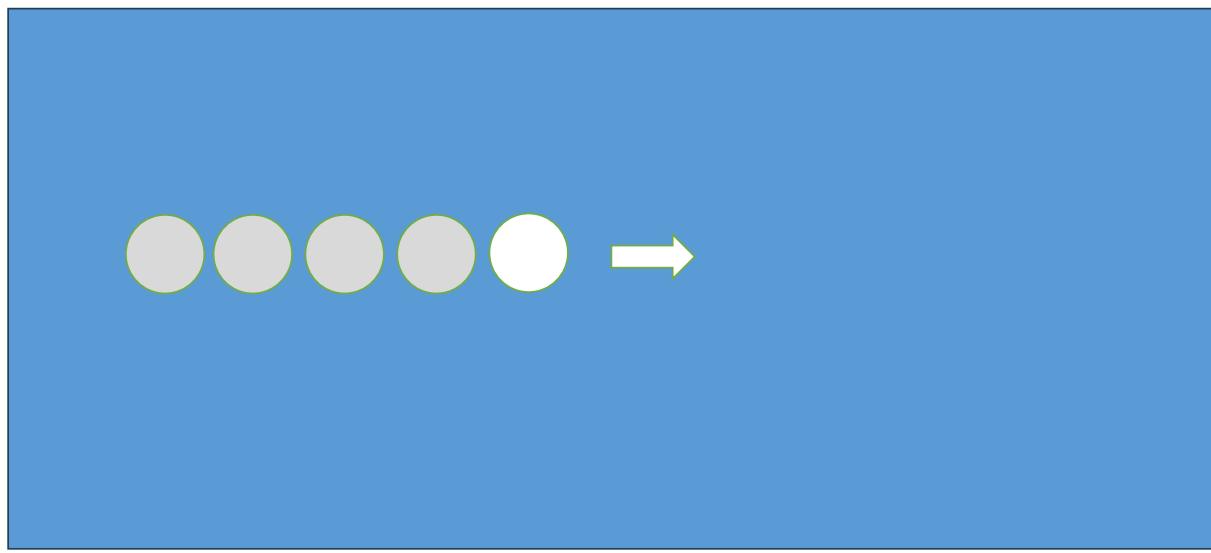
Given an image of a ball,  
can you predict where it will go next?



Source: MIT Introduction to Deep Learning 6.S191

4

Now, given a history of the ball's positions



Source: MIT Introduction to Deep Learning 6.S191

5

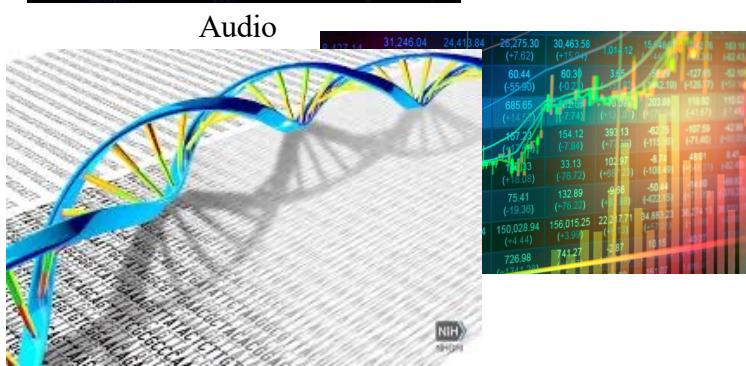
# Sequences & Sequence modelling



Theo had been staring at the same sentence so long the letters bled together.

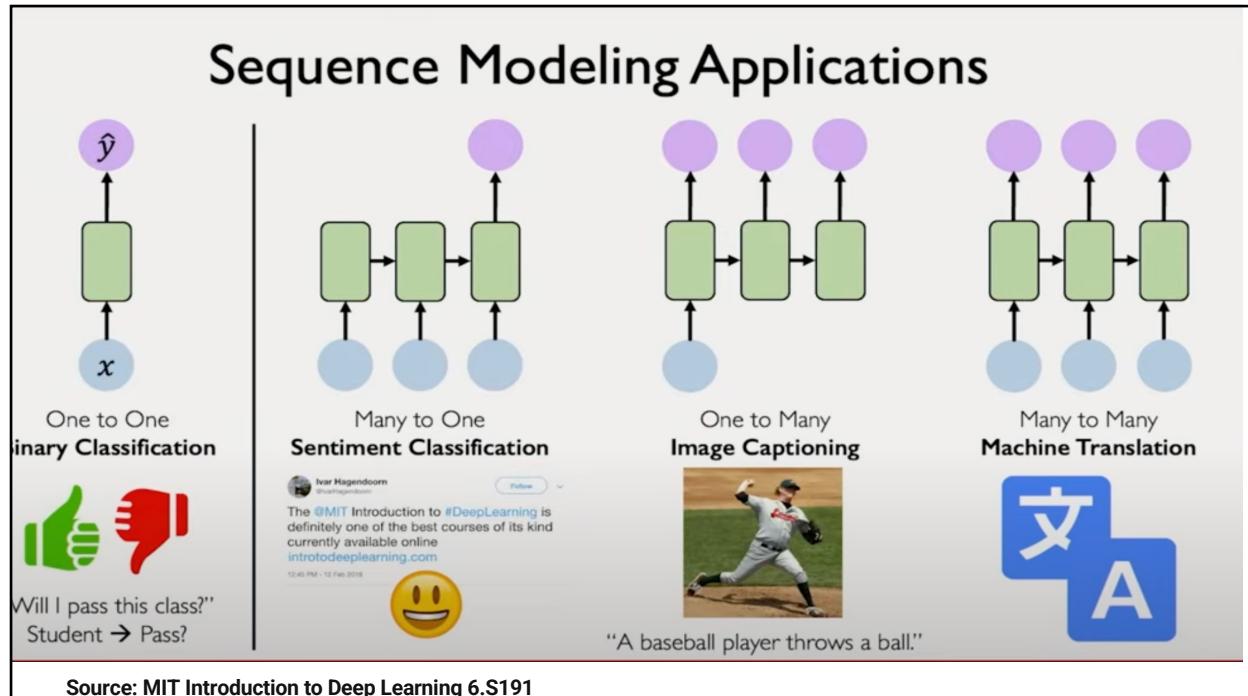
(*A Doorway on Lenox* by Ovett Chapman)

Text

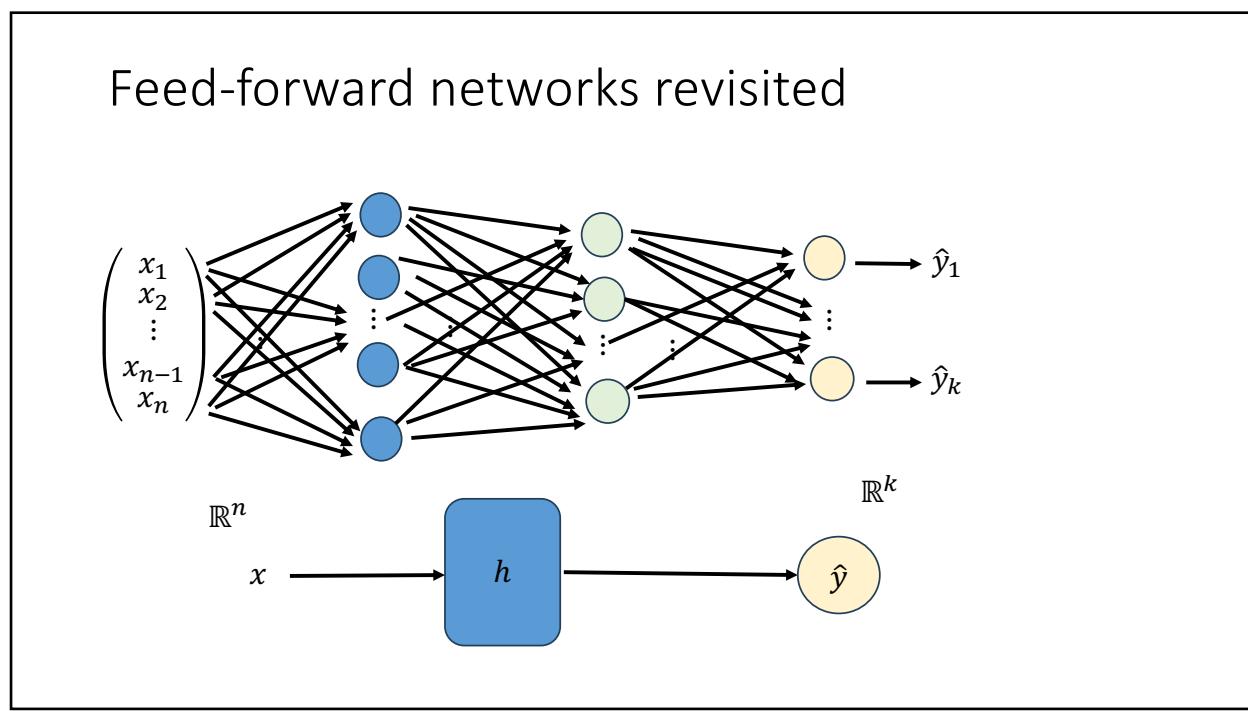


## Video

6

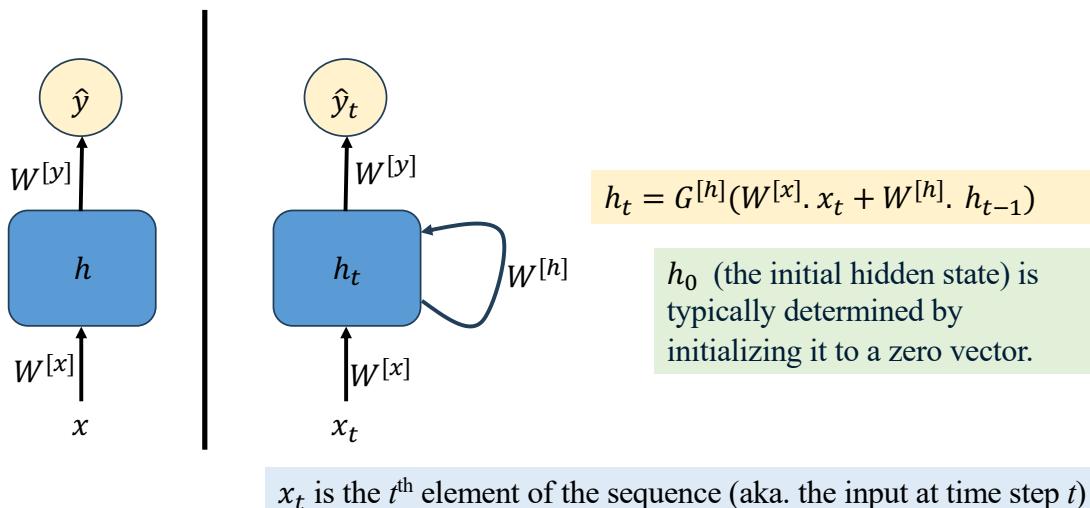


7



8

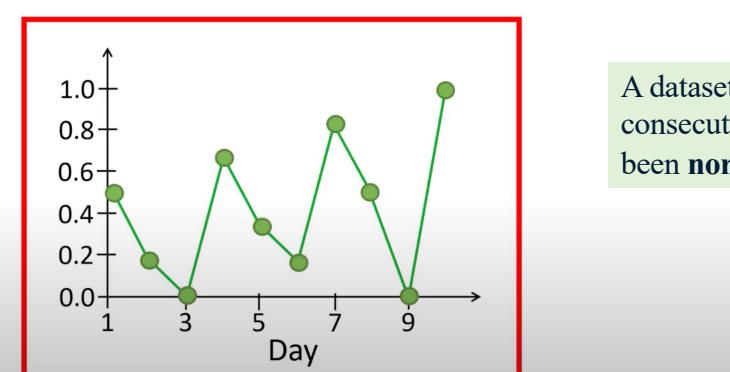
Extending traditional FF networks to deal with sequences



9

## Intuition of RNNs

	Mon	Tue	Wed	Thu	Fri	Mon	Tue	Wed	Thu	Fri
X	0.5	0.17	0.0	0.67	0.33	0.17	0.83	0.5	0.0	1.0



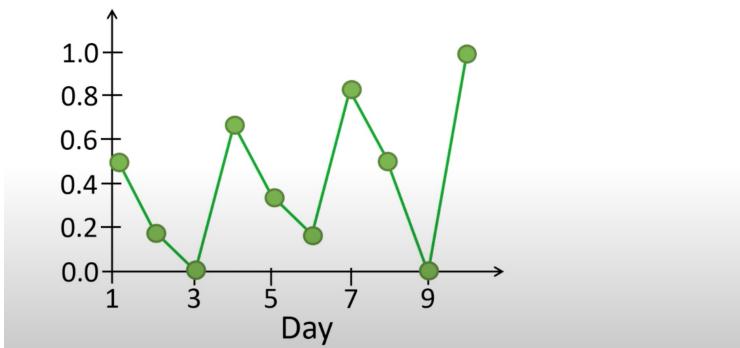
A dataset for stock prices (for 10 consecutive business days) that have been **normalized** into the  $[0,1]$  range.

Source: YouTube TileStats

10

## Intuition of RNNs

	Mon	Tue	Wed	Thu	Fri	Mon	Tue	Wed	Thu	Fri
X	0.5	0.17	0.0	0.67	0.33	0.17	0.83	0.5	0.0	1.0

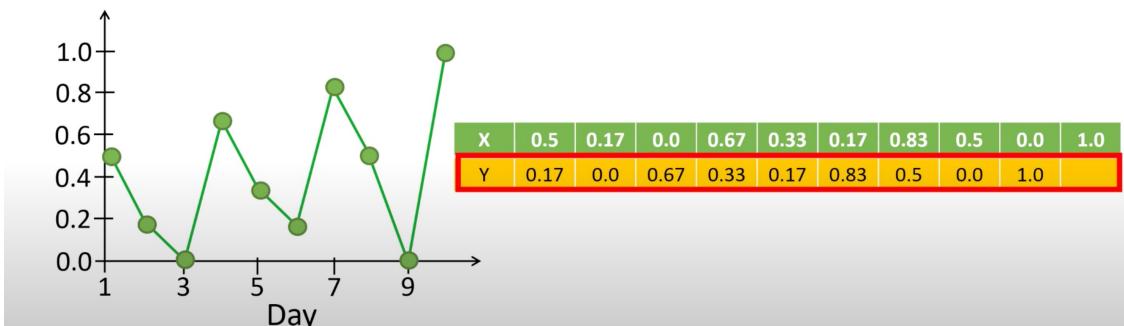


Source: YouTube TileStats

11

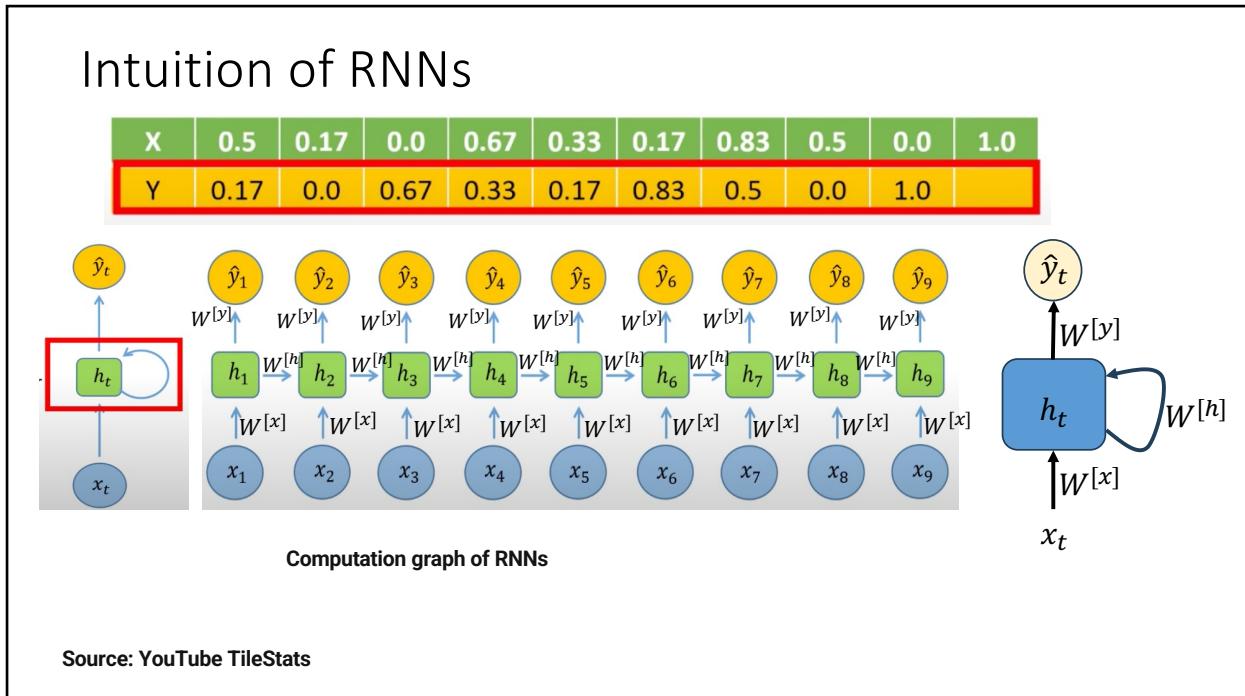
## Intuition of RNNs

	Mon	Tue	Wed	Thu	Fri	Mon	Tue	Wed	Thu	Fri
X	0.5	0.17	0.0	0.67	0.33	0.17	0.83	0.5	0.0	1.0
Y		0.17	0.0	0.67	0.33	0.17	0.83	0.5	0.0	1.0

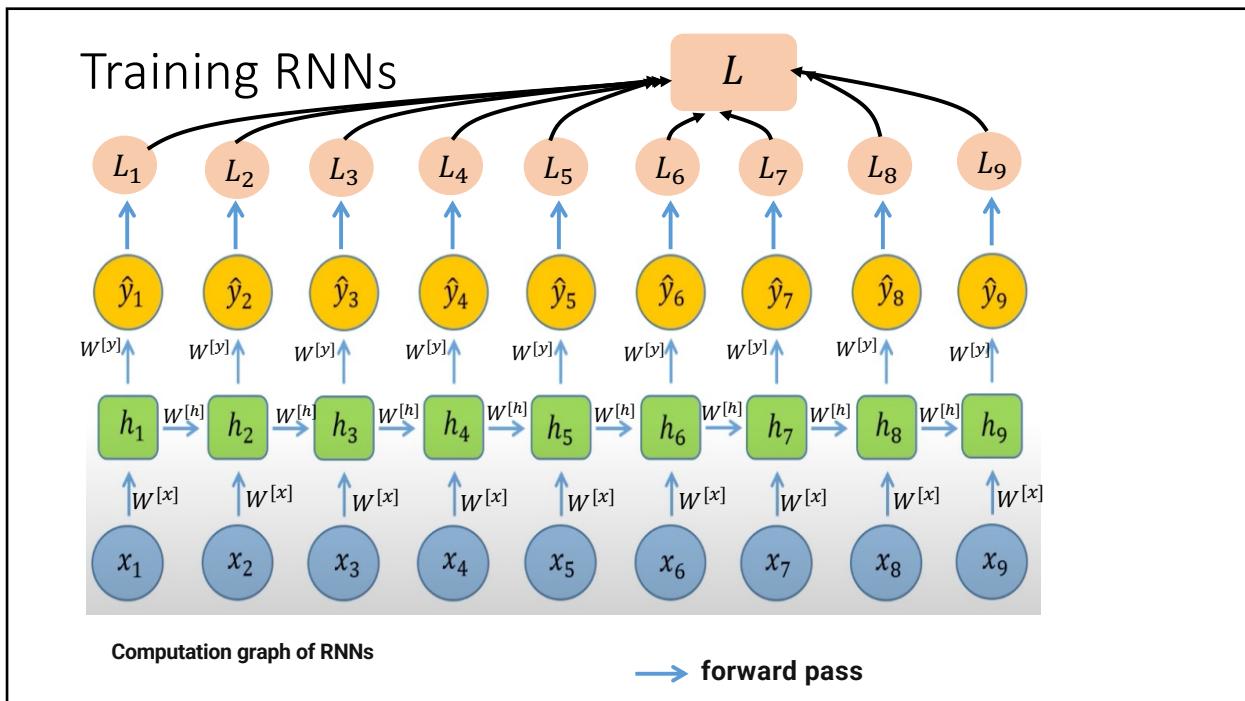


Source: YouTube TileStats

12

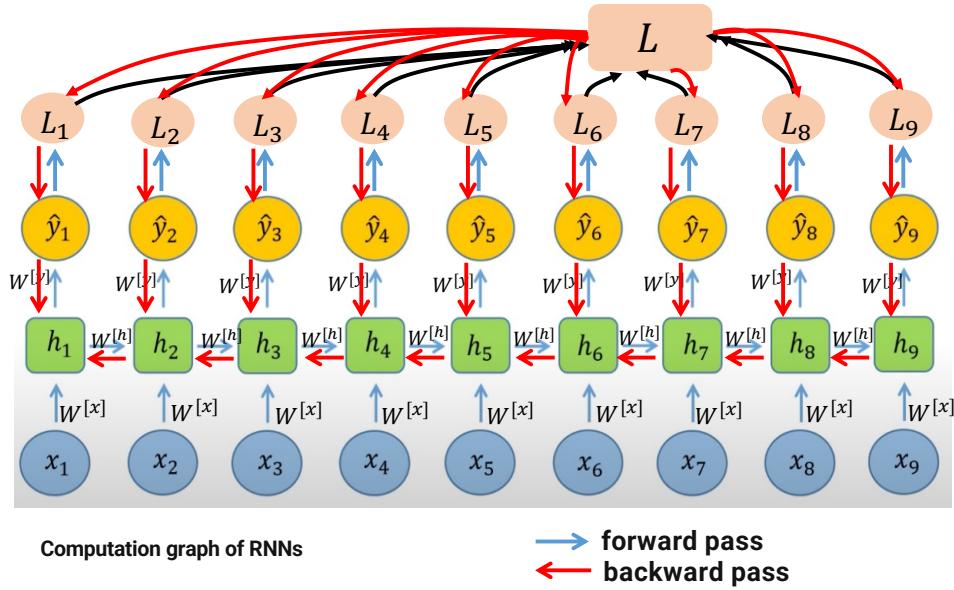


13



14

## RNNs: Backpropagation through Time (BPTT)



15

## RNNs: Backpropagation through Time (BPTT)

The loss function:

$$L = \frac{1}{T} \sum_{t=1}^T L_t = \frac{1}{T} \sum_{t=1}^T \text{loss}(y_t, \hat{y}_t)$$

We'll show the hardest derivative:

$$\frac{\partial L}{\partial W^{[h]}} = \frac{1}{T} \sum_{t=1}^T \frac{\partial \text{loss}(y_t, \hat{y}_t)}{\partial W^{[h]}} = \frac{1}{T} \sum_{t=1}^T \frac{\partial \text{loss}(y_t, \hat{y}_t)}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial h_t} \cdot \frac{\partial h_t}{\partial W^{[h]}}$$

The this effect o  
Very long recursive calculations of  $\frac{\partial h_t}{\partial W^{[h]}}$  for all  $t = 1, \dots, T$  ute the  
 $(-1))$

$$\frac{\partial h_t}{\partial W^{[h]}} = \frac{\partial f(W^{[h]}, W^{[x]}, x_t, h_{t-1})}{\partial W^{[h]}} + \frac{\partial f(W^{[h]}, W^{[x]}, x_t, h_{t-1})}{\partial h_{t-1}} \cdot \frac{\partial h_{t-1}}{\partial W^{[h]}}$$

16

## RNNs: Issues and Challenges

- Intensive computation (esp. for long sequences)
  - Truncation: Approximating the computation of  $\frac{\partial h_t}{\partial W^{[h]}}$  by going back only  $\tau$  steps.
- Numerical instability:
  - High powers of matrices can lead to divergent or vanishing eigenvalues.
  - exploding or vanishing gradients

**Many values > 1:**  
Exploding gradients  
→ Gradient clipping

**Many values < 1:**  
Vanishing gradients

- Activation function
- Parameter initialization
- Network architecture

17

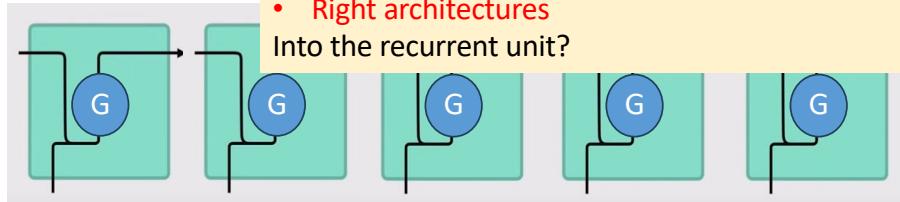
## RNNs: Limitations

- Encoding bottleneck (not all data are numerical, e.g., text, audio, video, etc.)
- Slow computation (due to no parallelization)
- No long memory

**Question:** Can we overcome (some of) these limitations by introducing the:

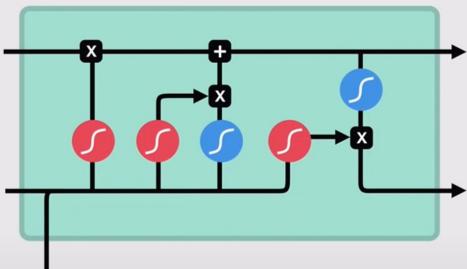
- Right activation functions
- Right architectures

Into the recurrent unit?

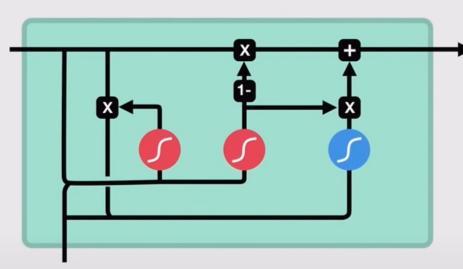


18

## LSTM & GRU



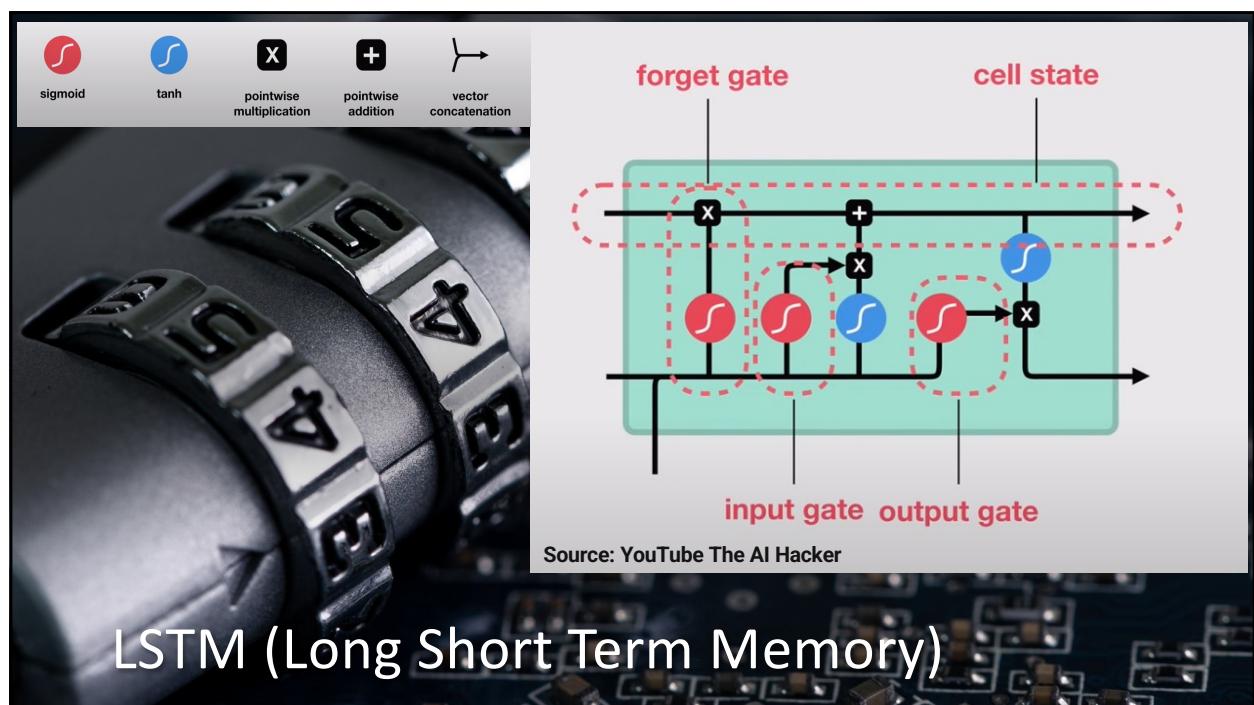
LSTM



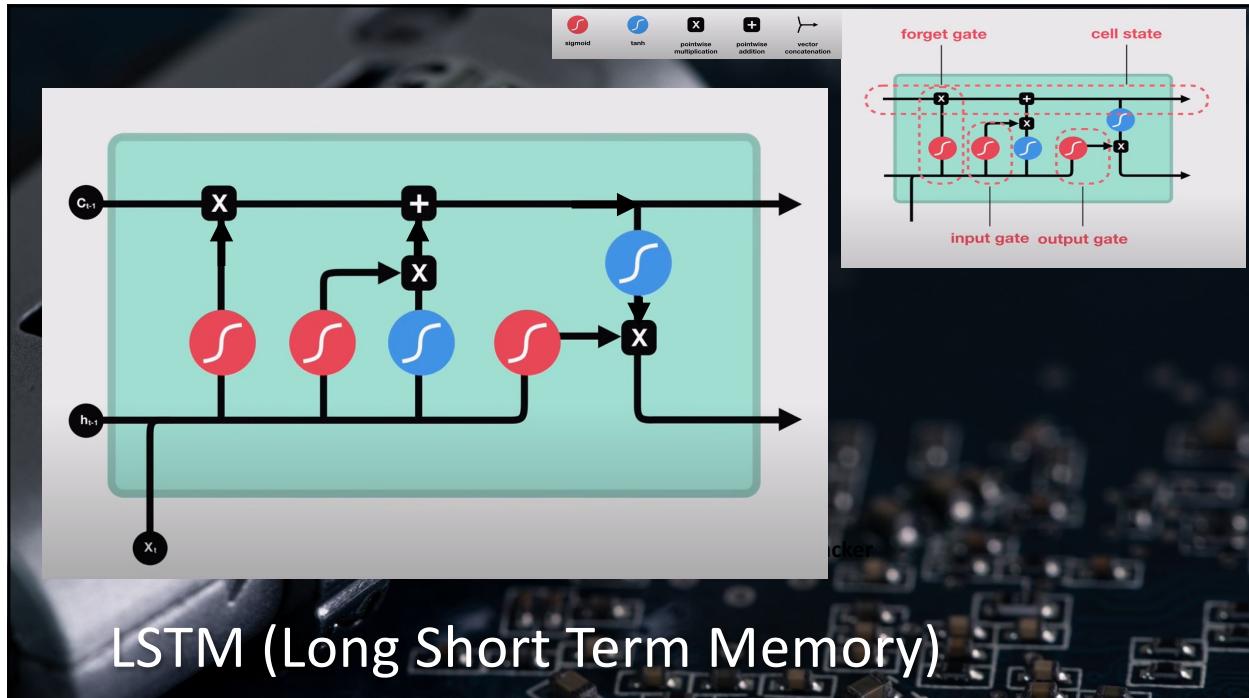
GRU

Source: YouTube The AI Hacker

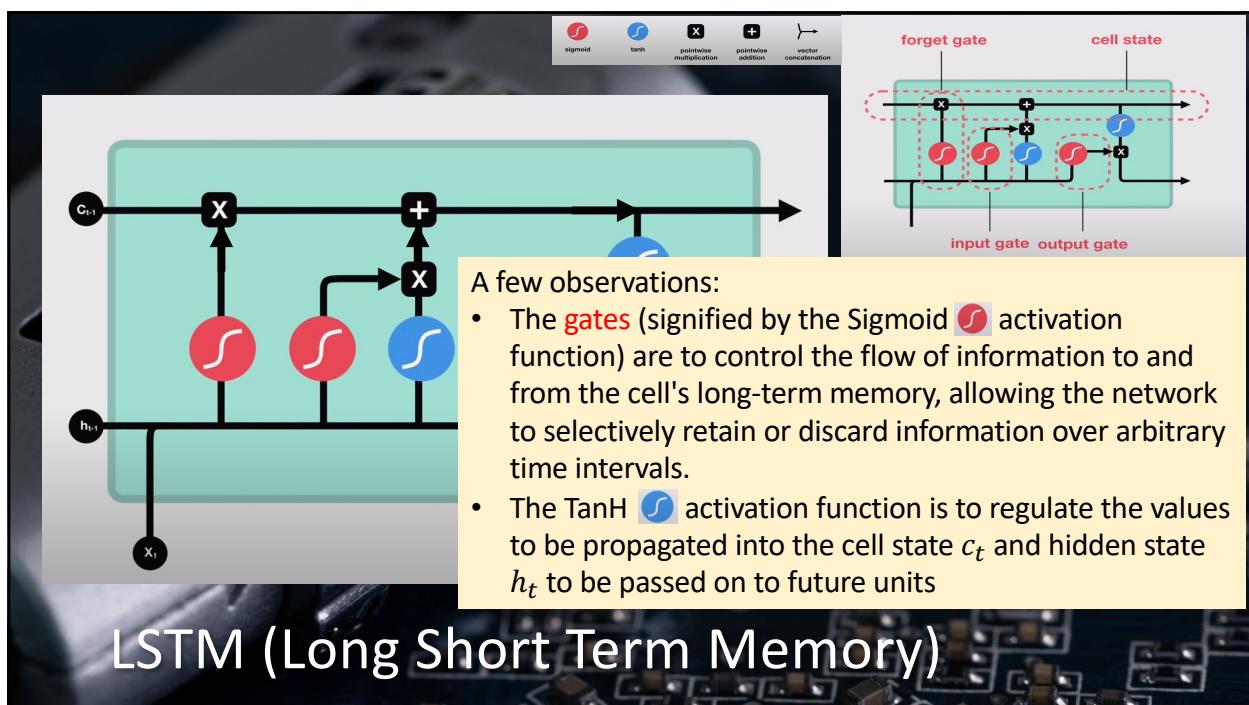
19



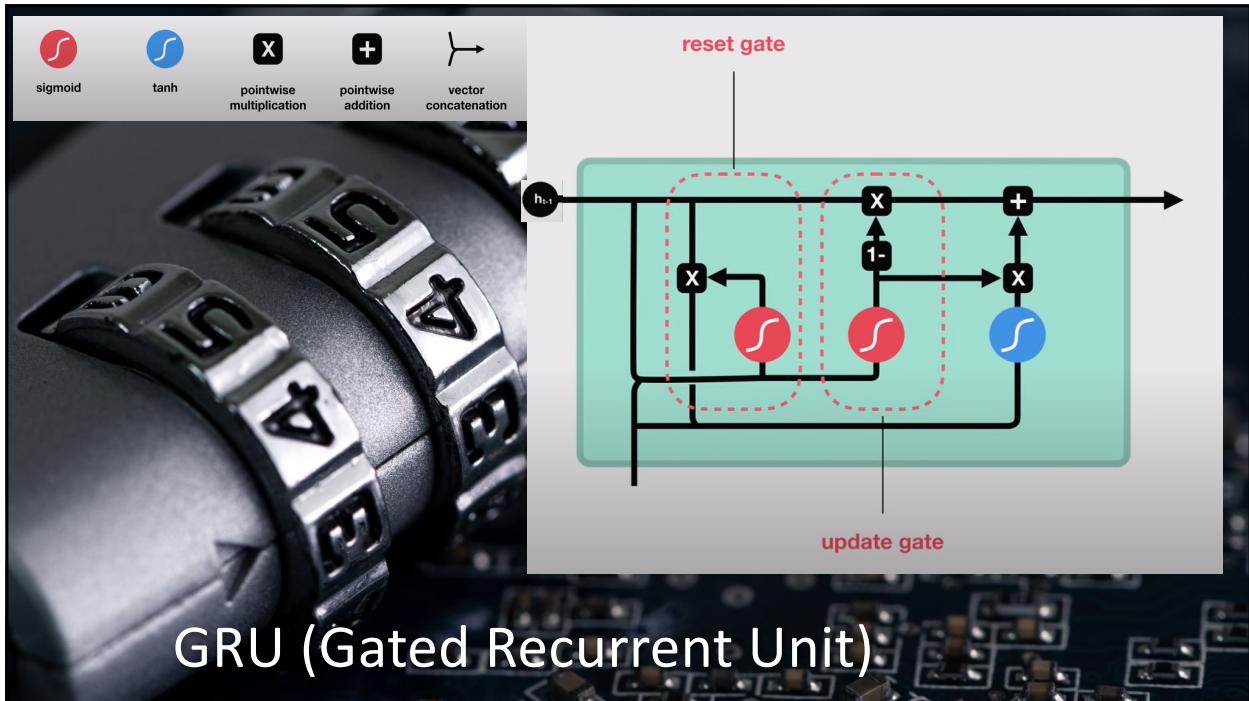
20



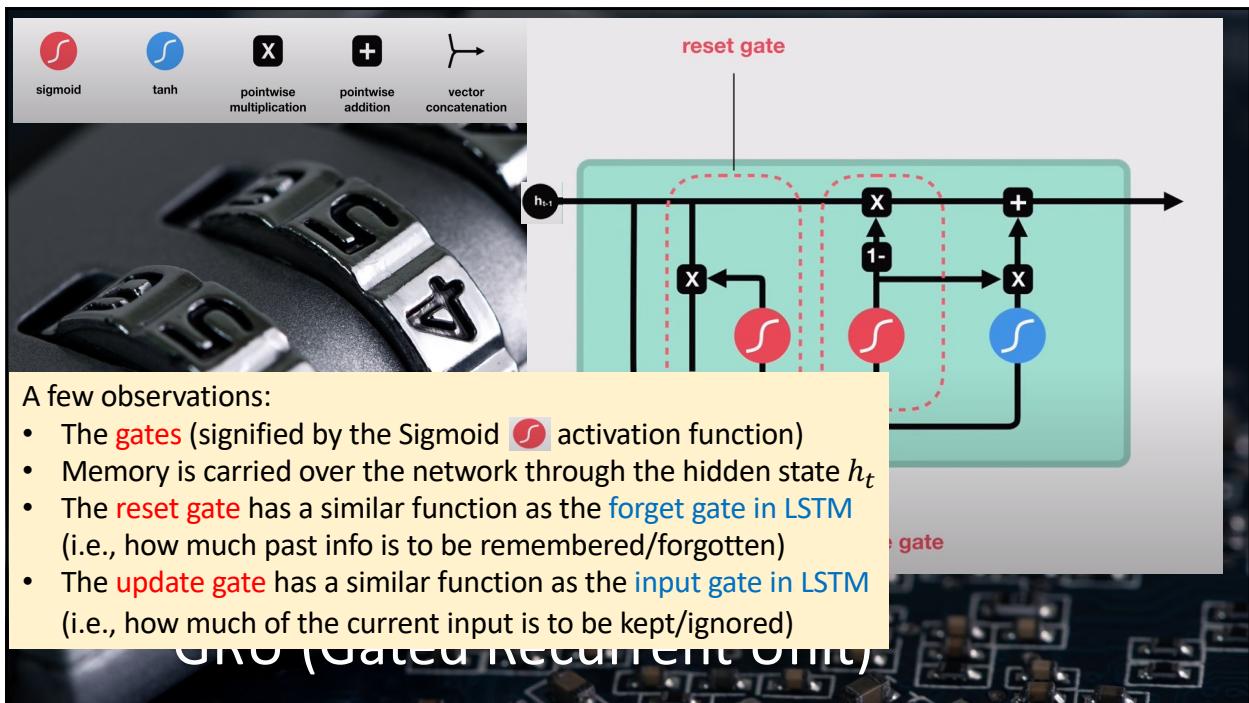
21



22



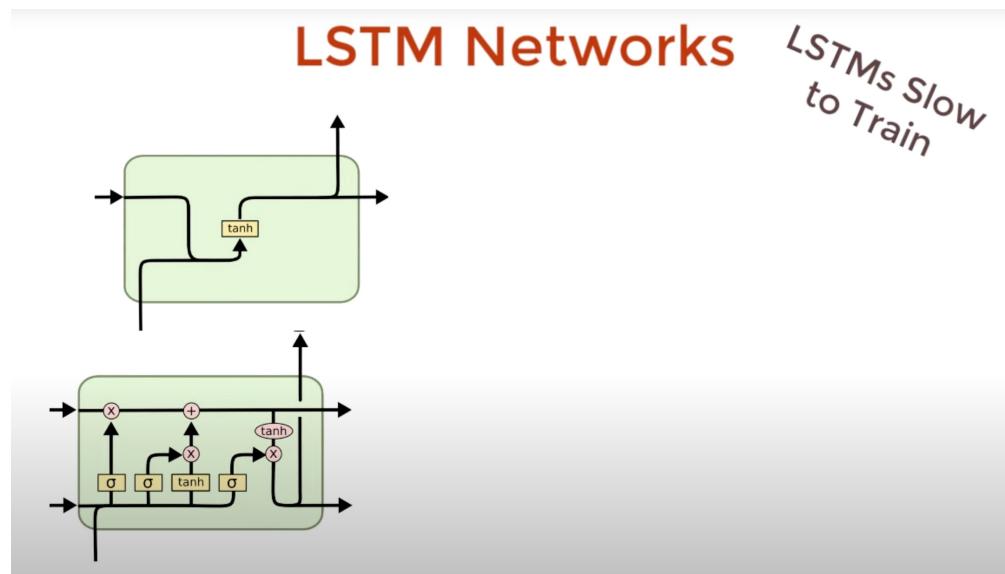
23



24

## LSTM drawbacks

Source: YouTube CodeEmporium



25

## LSTM drawbacks

Source: YouTube CodeEmporium

# Can we parallelize sequential data?

26

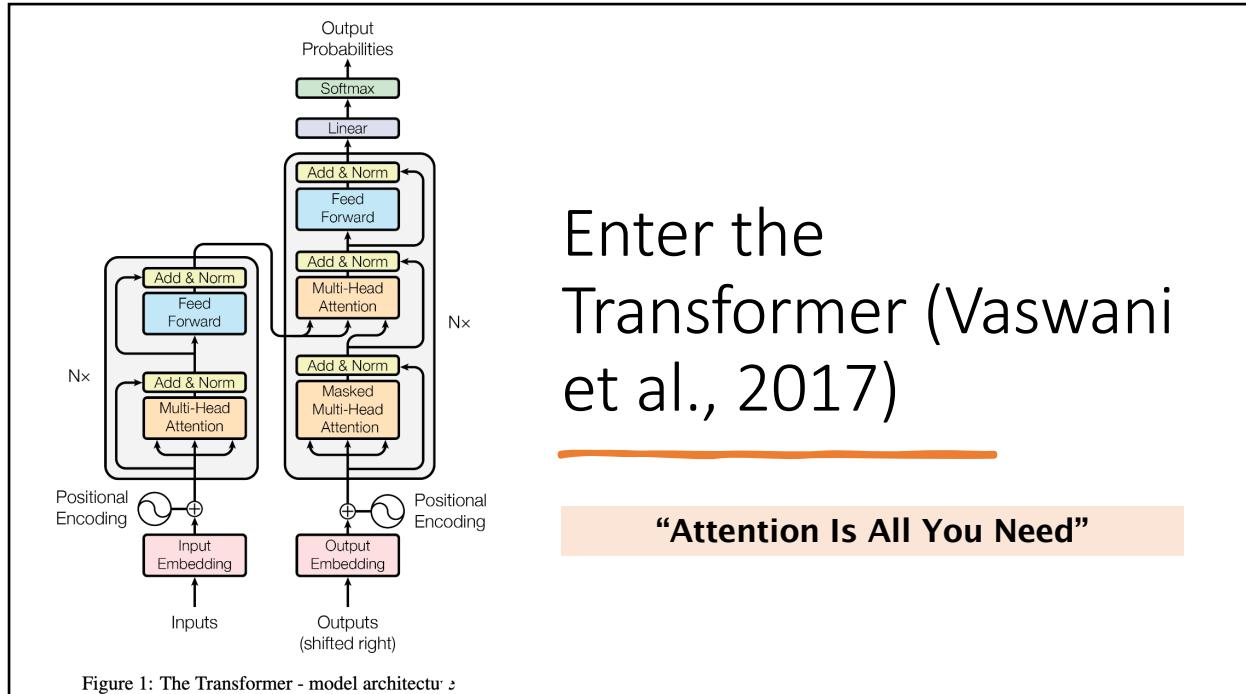
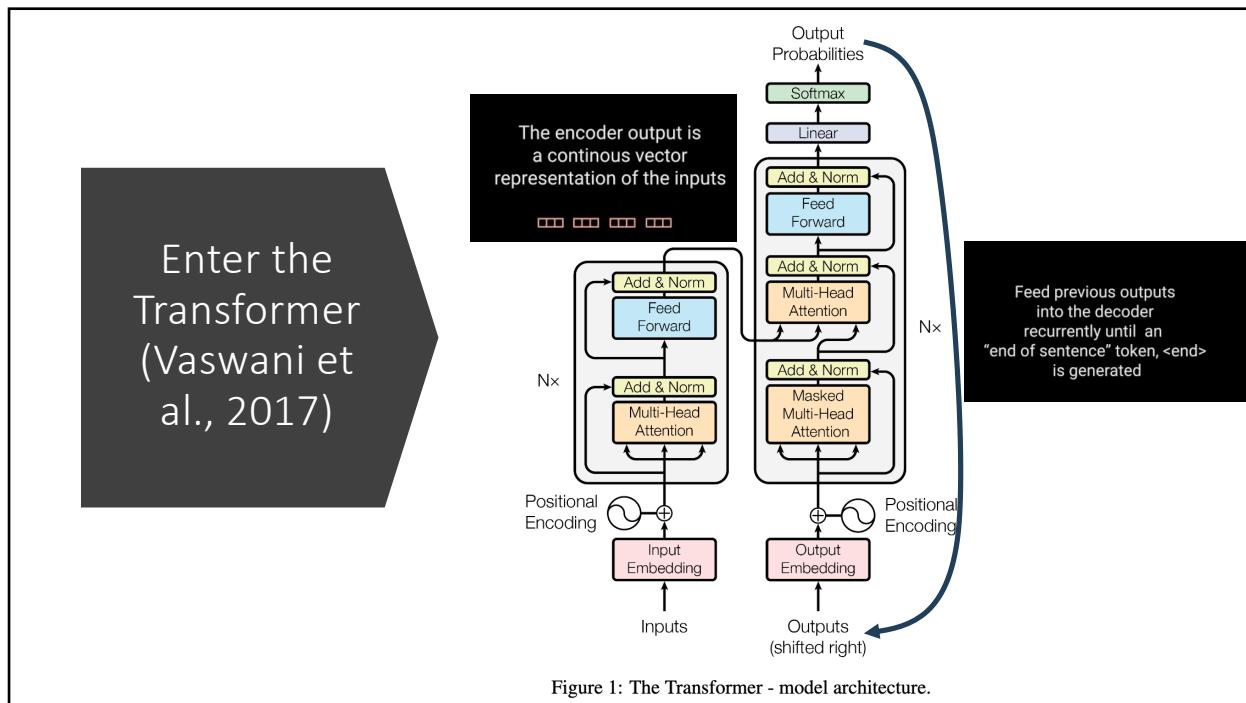


Figure 1: The Transformer - model architectur ↴

27

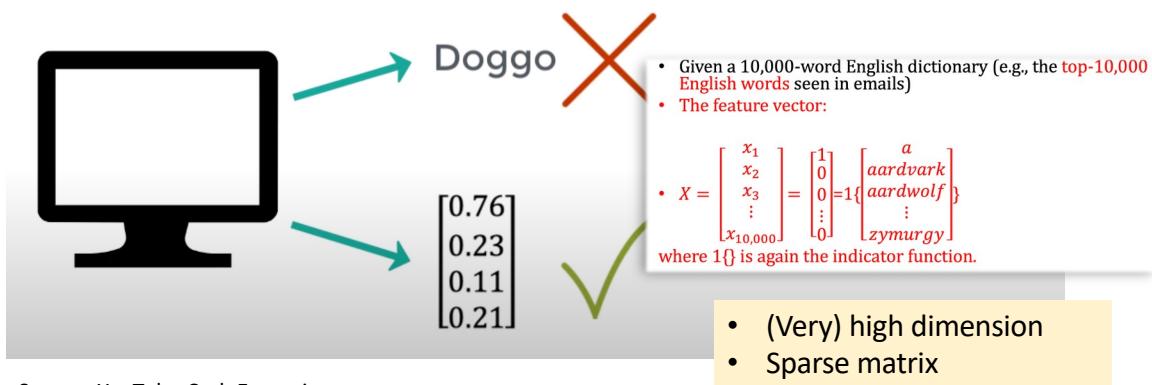


28

## Transformer Components

### Input Embedding

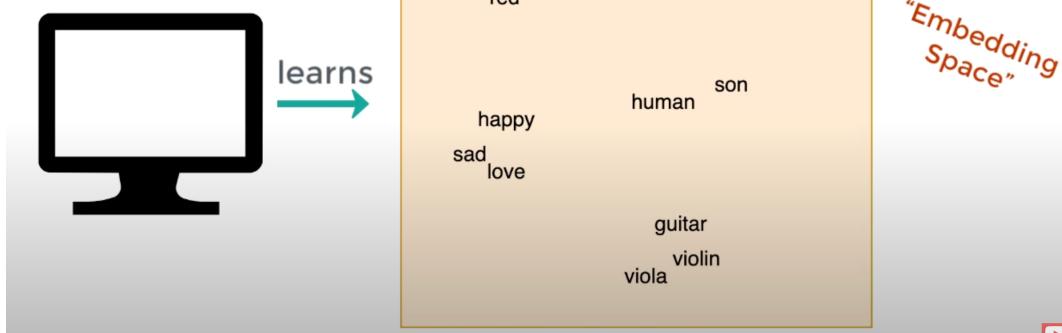
Remember our one-hot encoding in Week 5?



29

## Transformer Components

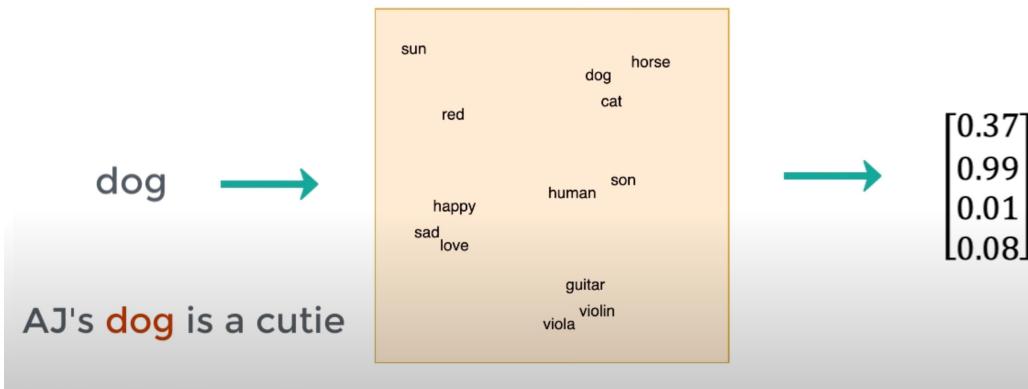
### Input Embedding



30

## Transformer Components

### Input Embedding



Source: YouTube CodeEmporium

31

## 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

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

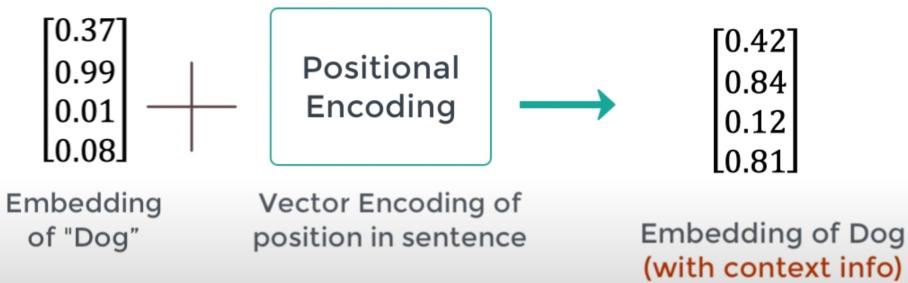
$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

Source: YouTube CodeEmporium

32

## Transformer Components

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

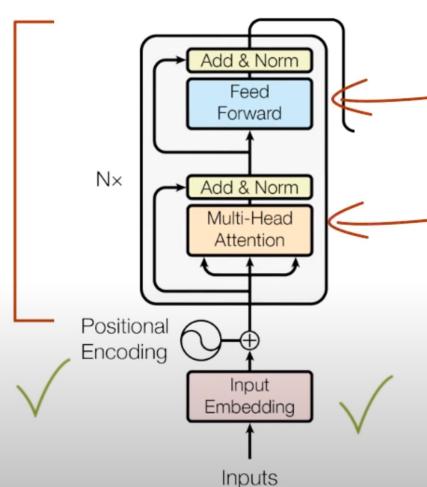


Source: YouTube CodeEmporium

33

## Transformer Components

Encoder Block

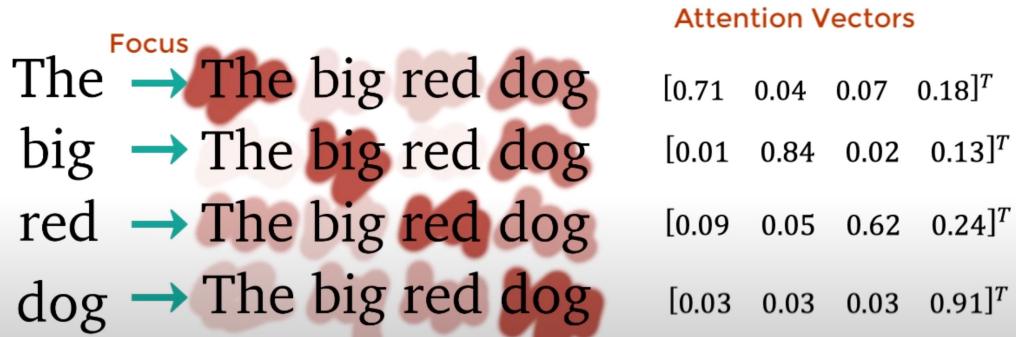


Source: YouTube CodeEmporium

34

## Transformer Components

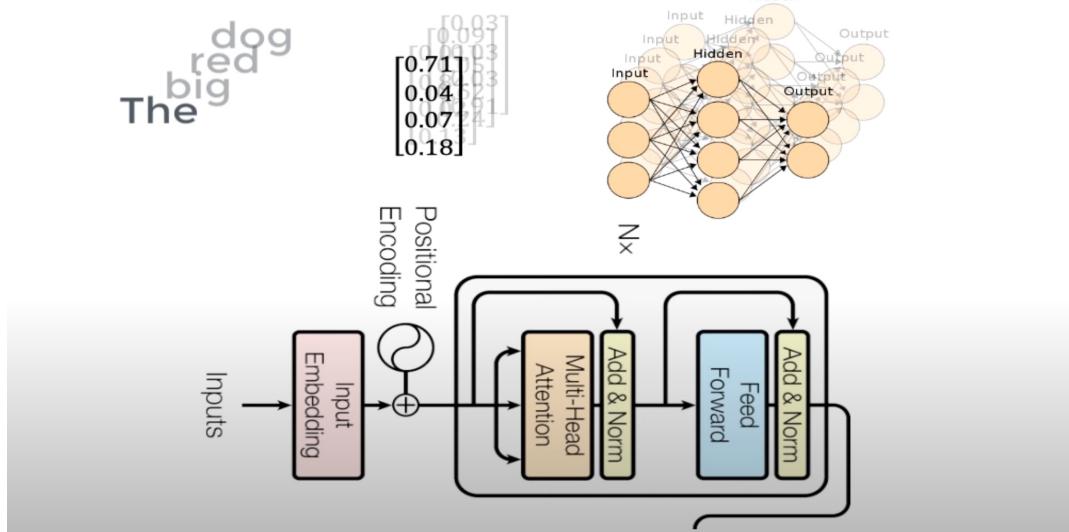
Attention : What part of the input should we focus?



Source: YouTube CodeEmporium

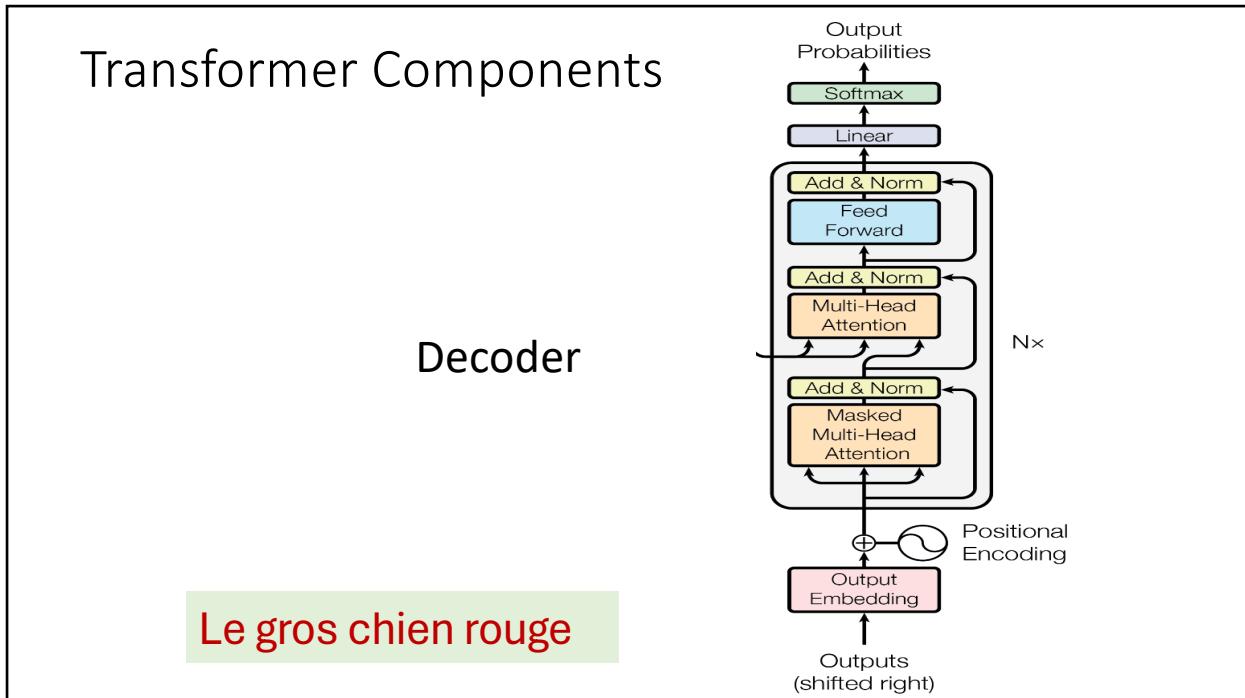
35

## Transformer Components

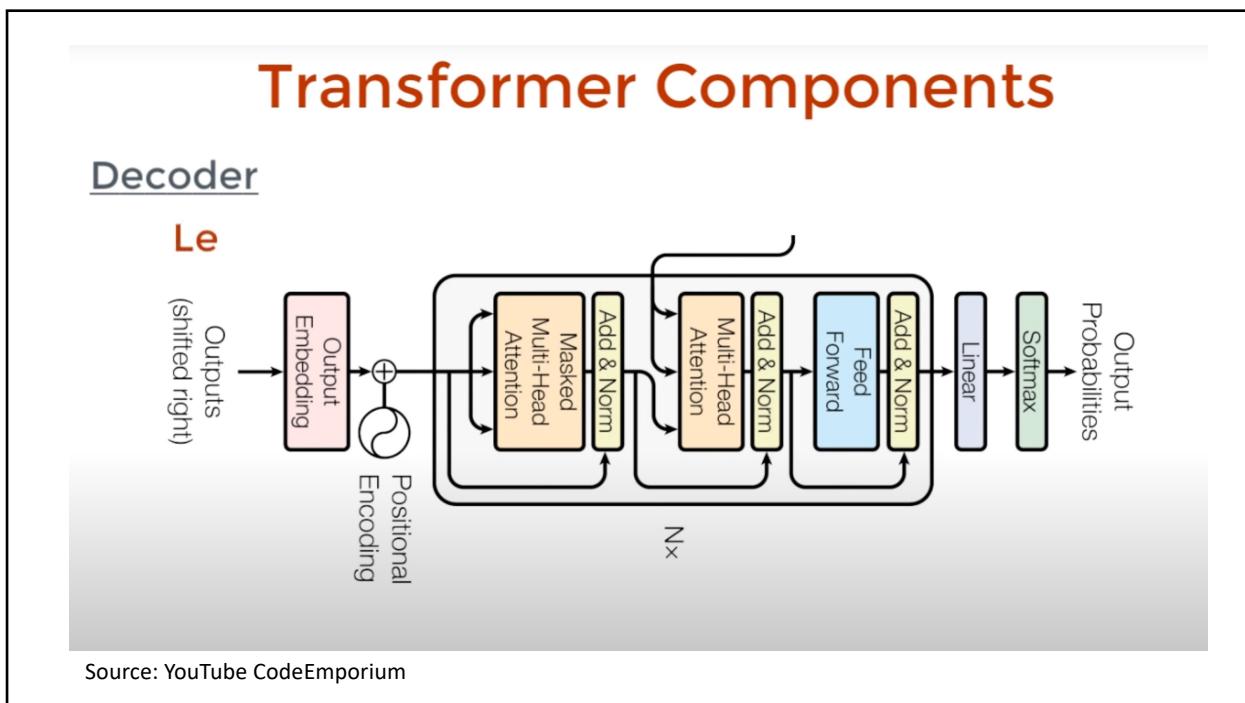


Source: YouTube CodeEmporium

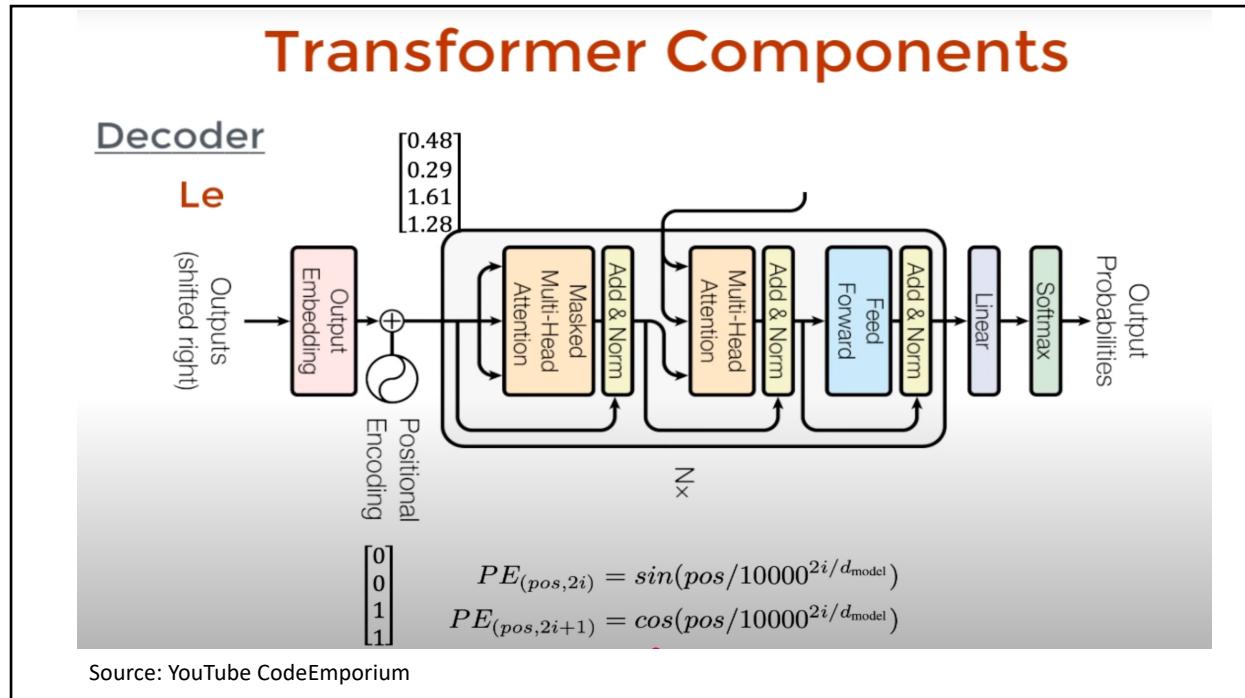
36



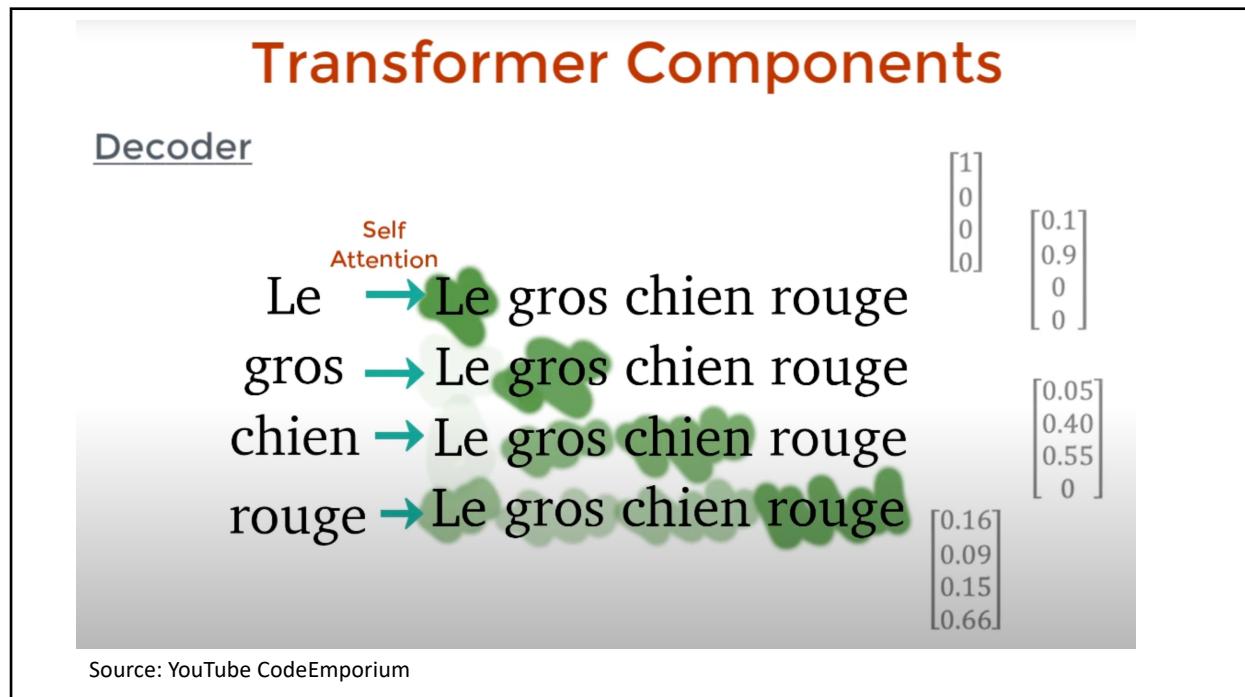
37



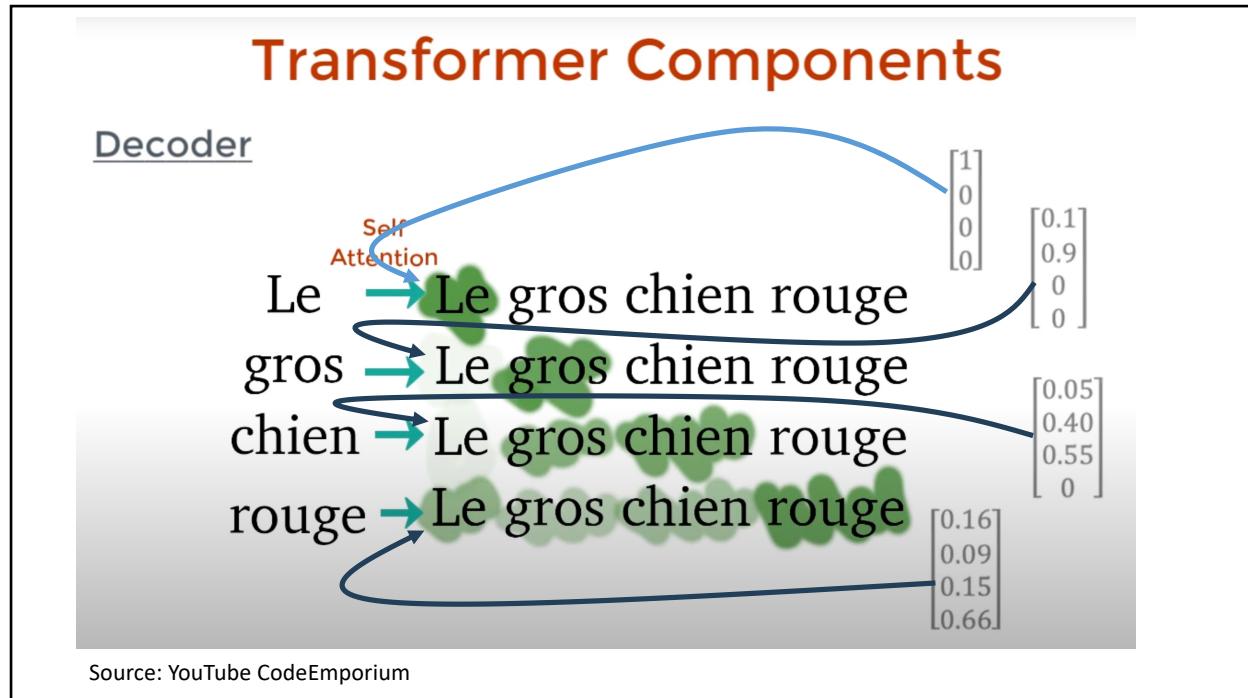
38



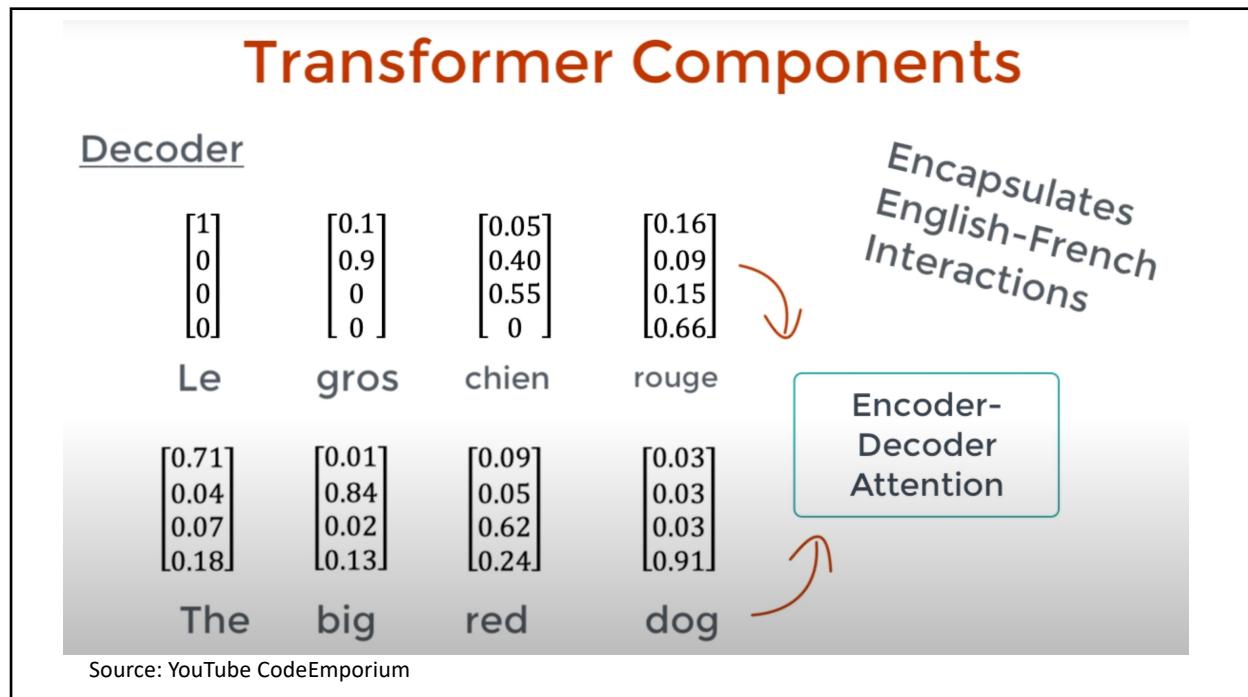
39



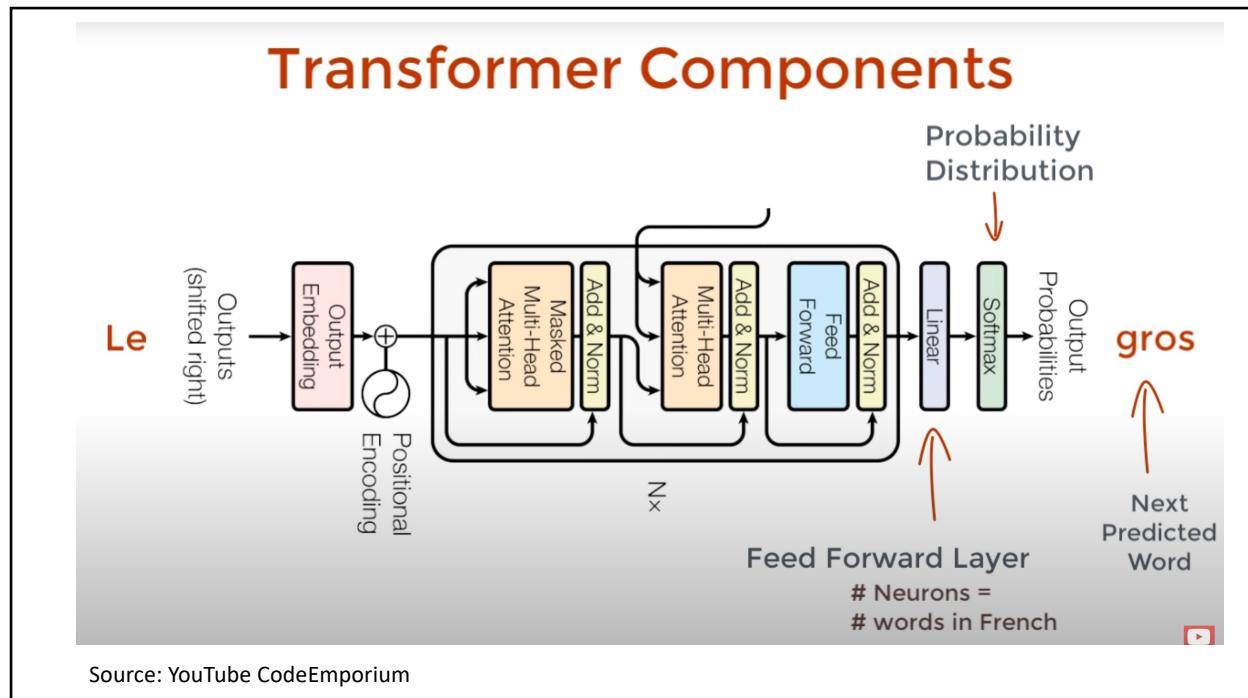
40



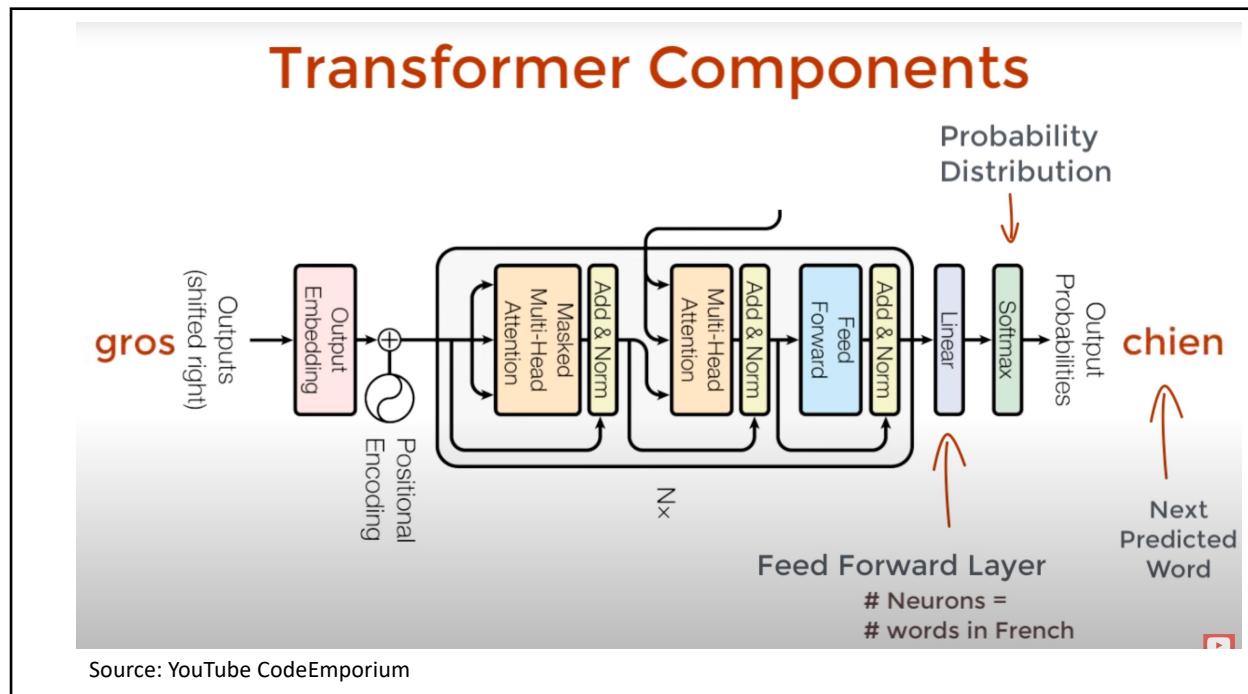
41



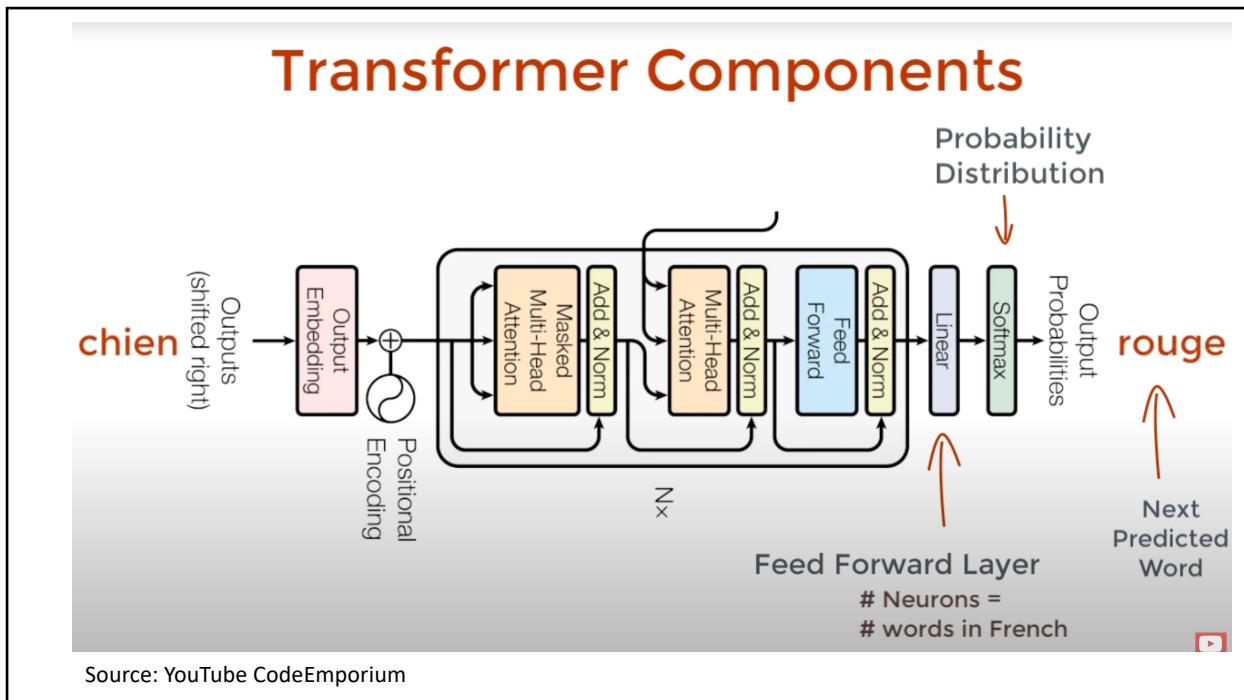
42



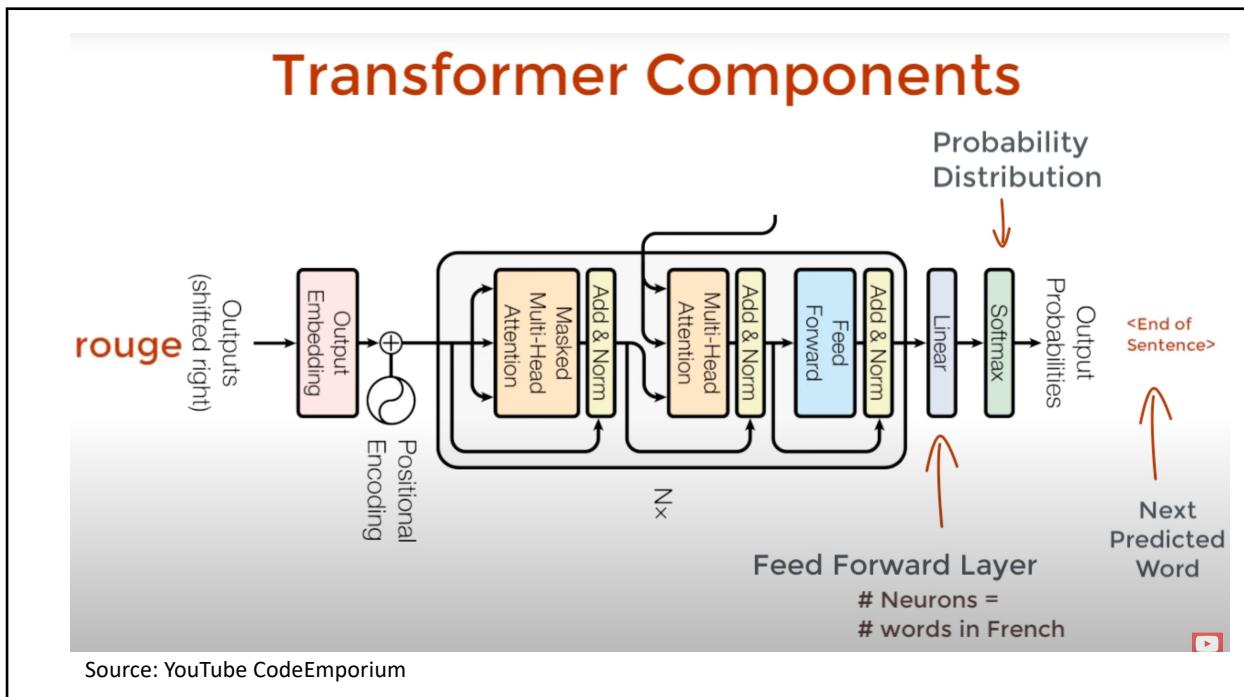
43



44



45



46

## Summary

- DL architectures for sequence data
- RNN
  - Encoding bottleneck
  - Slow computation (due to no parallelization)
  - No long memory
- LSTM & GRU
  - Advanced architectures for RNN
  - Address the issues of exploding/vanishing gradients
- Transformer
  - Foundation of most advanced DL systems
  - Architectural foundation for LLMs/LFMs
  - The combination of self-supervised learning (SSL) and transformers allows models to learn from large amounts of unlabeled data.