# Deep Learning Methods for NLP

**Machine Learning for Natural Language Processing**

**Pejman Rasti**

# Today Lecture Outline

- Deep Learning Framework

- The Multi-Layer Perceptron

- Recurrent Neural Network

- Self-Attention Mechanism and the Transformer Architecture

# Motivations

So far, we have seen, **techniques to represent tokens with vectors**

Given a certain representations of tokens:
➔ **How can we model a sequence of tokens to perform a specific task?**

In the past 10 years, a "new" class of machine learning techniques has become very popular and successful in NLP: **Deep Learning**

*In this session, we introduce Deep Learning with a focus[3] on the methods used in NLP*

# Framework

We want to model $(X_1, .., X_T)$ i.e. find the correct label $Y$

$$dnn_\theta : \qquad \mathbb{R}^{d,T} \qquad \rightarrow \mathbb{R}^p \ or \ [|0, K|]^p$$

$$(X_1, .., X_T) \mapsto \hat{Y}$$

- Output space is $\mathbb{R}^p$ for **Regression** tasks
- Output space is $[|0, K|]^p$ for **Classification** tasks

# Framework

We want to model $(X_1, .., X_T)$ i.e. find the correct label $Y$

$$dnn_\theta : \qquad \mathbb{R}^{d,T} \qquad \rightarrow \mathbb{R}^p \text{ or } [|0, K|]^p$$

$$(X_1, .., X_T) \mapsto \hat{Y}$$

**Questions: when we do Deep Learning…**

- **How do we define $dnn_\theta$ ?**
- **How do we train $dnn_\theta$ with data ?**

# Framework

Given a sequence of vectors $(X_1, .., X_T)$ we want to predict $Y$

$$dnn_\theta : \qquad \mathbb{R}^{d,T} \qquad \rightarrow \mathbb{R}^p \ or \ [\|0, K\|]^p$$

$$(X_1, .., X_T) \mapsto \hat{Y}$$

Most Deep Learning Models (all the ones we will use in this course):
- are **parametric**
- defined as a **composition of "simple" functions (linear & non-linear)**
- are trained in an **end-to-end** fashion with **backpropagation**

NB: In Deep Learning, **the parametrization** is called **the Architecture**

6

# Different Types of Architecture

**How can we define** our predictive function f ? $dnn_\theta$ ?

➔ Multi-Layer Perceptron

➔ Recurrent Layers

➔ Attention Layers

➔ Self-Attention Layers (in a Transformer Architecture)

# Different Types of Architecture

How can we define our predictive function f ? $dnn_\theta$ ?
➔ Multi-Layer Perceptron
➔ Recurrent Layers
➔ Attention Layers
➔ Self-Attention Layers (in a Transformer Architecture)

How do we **train our model**? (i.e. estimate the parameters of the model)
➔ **Stochastic Gradient Descent** also called **backpropagation** in this
  context

8

# The MultiLayer Perceptron (MLP)

*aka "the Most simple Deep Learning Architecture"*

The **MLP** works **on unidimensional data** (e.g. dimension $d$)

We present the **MLP in the regression case** (e.g. output space is $\mathbb{R}^2$ ))

$$dnn_\theta : \qquad \mathbb{R}^d \qquad \rightarrow \mathbb{R}^2$$
$$X \mapsto \hat{Y}$$

# The MultiLayer Perceptron (MLP)
*aka "the Most simple Deep Learning Architecture"*

The **MLP** works **on unidimensional data** (e.g. dimension $d$)
We present the **MLP in the regression case** (e.g. output space is $\mathbb{R}^2$ ))

$$dnn_{(W_1, b1, W_2, b_2)}(X) = W_2 \varphi_1(W_1 X + b_1) + b_2$$

$W_1, b_1, W_2$ and $b_2$ are trainable parameters. $W_1 \in \mathbb{R}^{\delta \times d}, b_1 \in \mathbb{R}^\delta, W_2 \in \mathbb{R}^{2 \times \delta}$ and $b_2 \in \mathbb{R}$

$\varphi_1$ is a fixed non-linear function, $\varphi_1 : \mathbb{R}^d \to \mathbb{R}^\delta$

# The MultiLayer Perceptron (MLP)

The **MLP** works **on unidimensional data** (e.g. dimension *d*)
We present the **MLP in the regression case** (e.g. output space is $\mathbb{R}^2$

$$dnn_{(W_1, b1, W_2, b_2)}(X) = W_2 \varphi_1(W_1 X + b_1) + b_2$$

$W_1, b_1, W_2$ and $b_2$ are trainable parameters. $W_1 \in \mathbb{R}^{\delta \times d}, b_1 \in \mathbb{R}^\delta, W_2 \in \mathbb{R}^{2 \times \delta}$ and $b_2 \in \mathbb{R}$

$\varphi_1$ is a fixed non-linear function, $\varphi_1 : \mathbb{R}^d \to \mathbb{R}^\delta$

➜ This model is a *2-layer* **MLP** model

# The MultiLayer Perceptron (MLP)

The **MLP** works **on unidimensional data** (e.g. dimension $d$)
We present the **MLP in the regression case** (e.g. output space is $\mathbb{R}^2$

$$dnn_{(W_1,b1,W_2,b_2)}(X) = W_2\varphi_1(W_1X + b_1) + b_2$$

$W_1, b_1, W_2$ and $b_2$ are trainable parameters. $W_1 \in \mathbb{R}^{\delta \times d}, b_1 \in \mathbb{R}^{\delta}, W_2 \in \mathbb{R}^{2 \times \delta}$ and $b_2 \in \mathbb{R}$

$\varphi_1$ is a fixed non-linear function, $\varphi_1 : \mathbb{R}^d \to \mathbb{R}^{\delta}$

➤ This model is a *2-layer* **MLP** model

➤ With **1 *hidden layer*** of $\delta$ dimension

# The MultiLayer Perceptron (MLP)

The **MLP** works **on unidimensional data** (e.g. dimension *d*)
We present the **MLP in the regression case** (e.g. output space is $\mathbb{R}^2$

$$dnn_{(W_1, b1, W_2, b_2)}(X) = W_2 \varphi_1 (W_1 X + b_1) + b_2$$

$W_1, b_1, W_2$ and $b_2$ are trainable parameters. $W_1 \in \mathbb{R}^{\delta \times d}, b_1 \in \mathbb{R}^\delta, W_2 \in \mathbb{R}^{2 \times \delta}$ and $b_2 \in \mathbb{R}$

$\varphi_1$ is a fixed non-linear function, $\varphi_1 : \mathbb{R}^d \to \mathbb{R}^\delta$

➜    This model is a *2-layer* **MLP** model
➜    With **1 *hidden layer*** of dimension $\delta$
➜    Taking as input a vector of **dimension *d*** to output a vector of **dimension 2**

# The MultiLayer Perceptron (MLP)

The **MLP** works **on unidimensional data** (e.g. dimension $d$)
We present the **MLP in the regression case** (e.g. output space is $\mathbb{R}^2$

$$dnn_{(W_1, b1, W_2, b_2)}(X) = W_2 \varphi_1(W_1 X + b_1) + b_2$$

$W_1, b_1, W_2$ and $b_2$ are trainable parameters. $W_1 \in \mathbb{R}^{\delta \times d}, b_1 \in \mathbb{R}^\delta, W_2 \in \mathbb{R}^{2 \times \delta}$ and $b_2 \in \mathbb{R}$

$\varphi_1$ is a fixed non-linear function, $\varphi_1 : \mathbb{R}^d \rightarrow \mathbb{R}^\delta$

➔ This model is a *2-layer* **MLP** model
➔ With **1 *hidden layer*** of dimensic $\delta$
➔ Taking as input a vector of **dimension $d$** to output a vector of **dimension 2**
➔ Such a model is also referred to as a *Feed-Forward* **Neural Network (FNN)**

# The MultiLayer Perceptron: Diagram View



Neurons

Figure from (R. Rezvani et. al. 2012)

In Deep Learning, it is usual to represent equations **with diagrams**

# The MultiLayer Perceptron: Diagram View

Output Y

Feed-Forward Layer

Input X

In Deep Learning, it is usual to represent equations **with diagrams**

# The MultiLayer Perceptron:

We have defined a 2-layers MLP model
We can define in the same way a **3-layers**, **4-layers**, **L-layers** MLP

$$dnn_{(W_i\ b_i,\ i \in [|1,L|])}(X) = W_L \varphi_{L-1}(...\varphi_2 \circ W_2 \varphi_1(W_1 X + b_1) + b_2)...) + b_L$$

$W_l$ and $b_l$ are trainable parameters. $W_l \in \mathbb{R}^{\delta_{l-1} \times \delta_l}$, $b_l \in \mathbb{R}^{\delta_l}$, with $\delta_l \in \mathbb{N}^*, \forall l \in [|1, L|]$

$\varphi_l$ fixed non-linear functions, $\varphi_l : \mathbb{R}^{\delta_{l-1}} \to \mathbb{R}^{\delta_l}, \forall l \in [|1, L-1|]$ <sup>17</sup>

# Recurrent Neural Network

# Example

Cooking Schedule

Monday   Tuesday   Wednesday   Thursday   Friday   Saturday

# Example

## Vectors


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


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


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

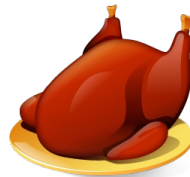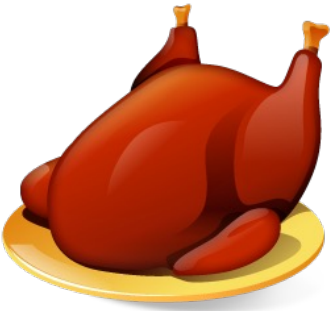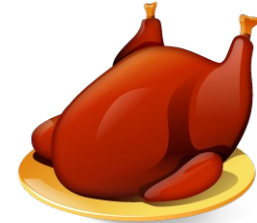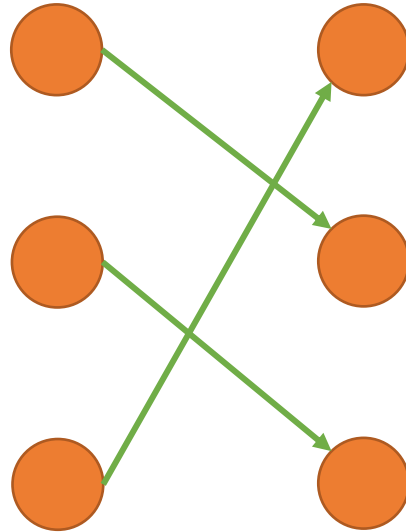# Example


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


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

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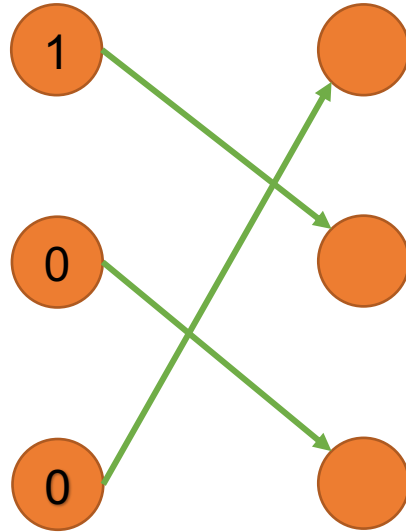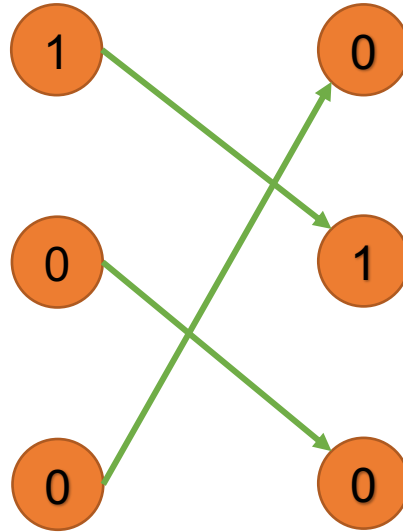

$$\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

# Example



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

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



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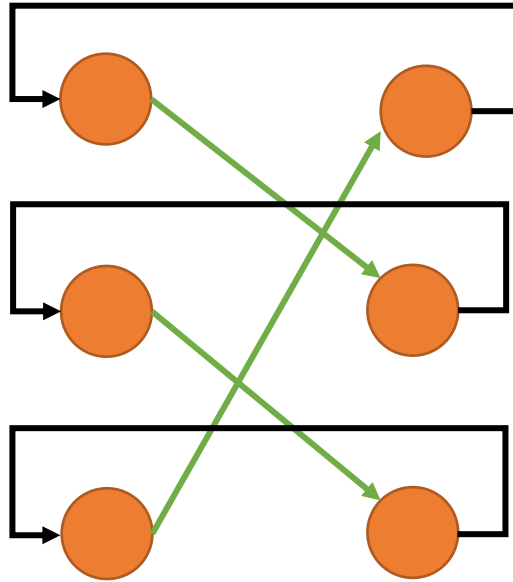


$$\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

# Example

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

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

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

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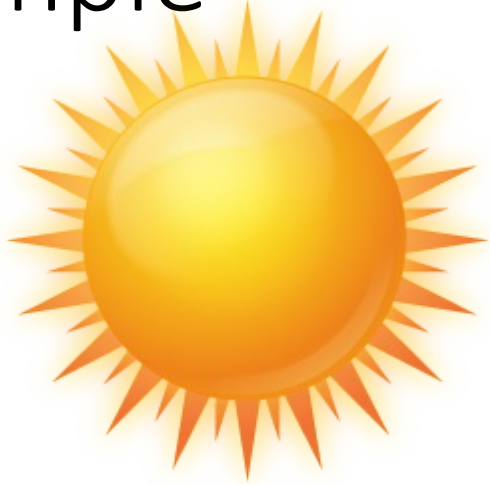
 $\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$ 

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

# Example

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

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

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

# Example

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



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

# Example

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



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

# Example

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



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



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

# Example

Simple (Recurrent) neural network

# Example



Weather

Same as yesterday

Next dish
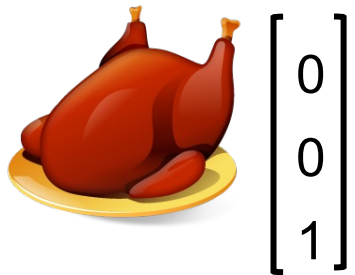
# Example

## Cooking Schedule

# Example

# Example

## Vectors

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

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

# Example

## More Complicated RNN

$$\begin{bmatrix} 1\ 0\ 0 \\ 0\ 1\ 0 \\ 0\ 0\ 1 \\ \hline 0\ 0\ 1 \\ 1\ 0\ 0 \\ 0\ 1\ 0 \end{bmatrix} \quad + \quad \begin{bmatrix} 1\ 0 \\ 1\ 0 \\ 1\ 0 \\ \hline 0\ 1 \\ 0\ 1 \\ 0\ 1 \end{bmatrix} \quad \Rightarrow$$

Food       Add       Weather       Merge

# Example



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

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

$$\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$
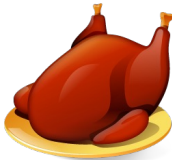
$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ \hline 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \hline 0 \\ 1 \\ 0 \end{bmatrix}$$
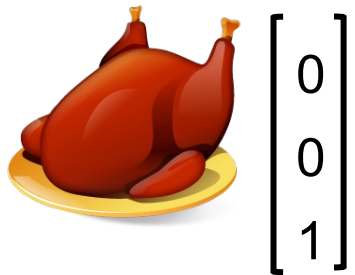
Food

Same

Next Day

# Example

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

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

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

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ \hline 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \hline 0 \\ 0 \\ 1 \end{bmatrix}$$
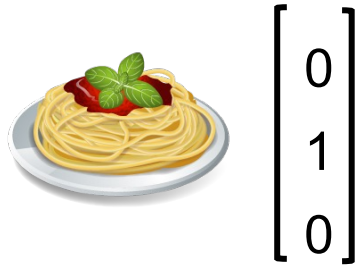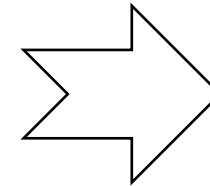
Food

Same

Next Day

# Example



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

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

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

$$\begin{bmatrix} \begin{array}{ccc} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ \hline 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{array} \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} \begin{array}{c} 0 \\ 0 \\ 1 \\ \hline 1 \\ 0 \\ 0 \end{array} \end{bmatrix}$$

Food

Same

Next Day

# Example

## More Complicated RNN

$$\begin{bmatrix} 1\ 0\ 0 \\ 0\ 1\ 0 \\ \underline{0\ 0\ 1} \\ 0\ 0\ 1 \\ 1\ 0\ 0 \\ 0\ 1\ 0 \end{bmatrix} \quad + \quad \begin{bmatrix} 1\ 0 \\ 1\ 0 \\ \underline{1\ 0} \\ 0\ 1 \\ 0\ 1 \\ 0\ 1 \end{bmatrix} \quad \Longrightarrow$$

Food         Add         Weather         Merge

# Example



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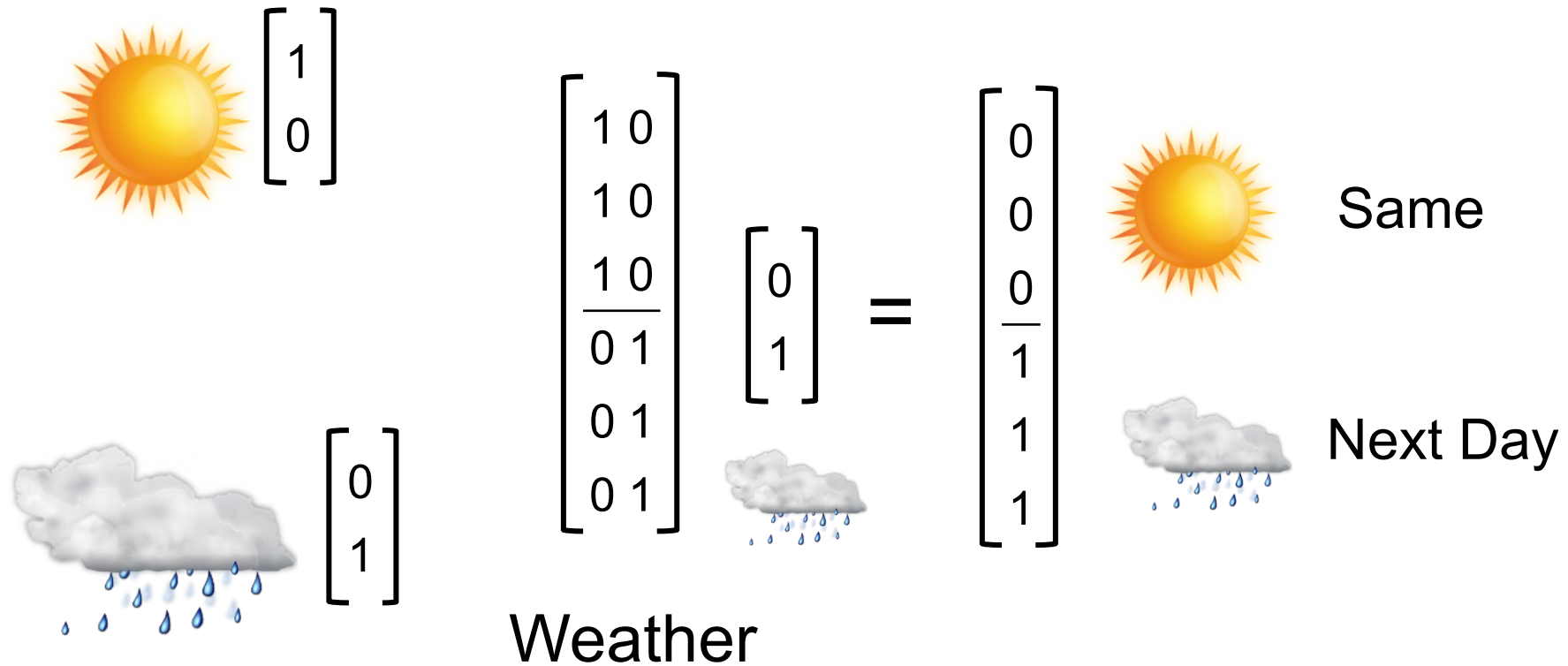
$$\begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ \hline 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ \hline 0 \\ 0 \\ 0 \end{bmatrix}$$

Weather

Same

Next Day

# Example



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

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

$$\begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ \hline 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \hline 1 \\ 1 \\ 1 \end{bmatrix}$$

Same

Next Day

Weather

# Example

## More Complicated RNN
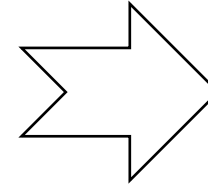
$$\begin{bmatrix} 1\ 0\ 0 \\ 0\ 1\ 0 \\ 0\ 0\ 1 \\ \hline 0\ 0\ 1 \\ 1\ 0\ 0 \\ 0\ 1\ 0 \end{bmatrix}$$

**+**

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

⇨

Food            Add            Weather        Merge

# Example

# Example

## More Complicated RNN

$$\begin{bmatrix} 1\ 0\ 0 \\ 0\ 1\ 0 \\ 0\ 0\ 1 \\ \hline 0\ 0\ 1 \\ 1\ 0\ 0 \\ 0\ 1\ 0 \end{bmatrix} \quad + \quad \begin{bmatrix} 1\ 0 \\ 1\ 0 \\ 1\ 0 \\ \hline 0\ 1 \\ 0\ 1 \\ 0\ 1 \end{bmatrix}$$

Food          Add          Weather          Merge

# Example

Merge

$$\begin{bmatrix} 1 \\ 0 \\ \hline 0 \\ \hline 1 \\ 2 \\ 1 \end{bmatrix}$$ Non-Linear $$\begin{bmatrix} 0 \\ 0 \\ \hline 0 \\ \hline 0 \\ 1 \\ 0 \end{bmatrix}$$ $$\begin{bmatrix} 0+0 \\ 0+1 \\ 0+0 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

# Example

## Recurrent neural network



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

Food

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

Weather

Add

Non-Linear

Merge

# Example



## Recurrent neural network

Food

Weather

Add

Non-Linear

Merge

# Example

## Recurrent neural network



Food

Weather

Add

Non-Linear

Merge

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

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

1
0
0
0
1
0

0
0
0
1
1
1

# Example

## Recurrent neural network



Food

Weather

Add

Non-Linear

Merge

48

# Example

## Recurrent neural network



Food

Weather

Add

Non-Linear

Merge

# Example

## Recurrent neural network

# Example

## Recurrent neural network



Food

Weather

Add

Non-Linear

Merge

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

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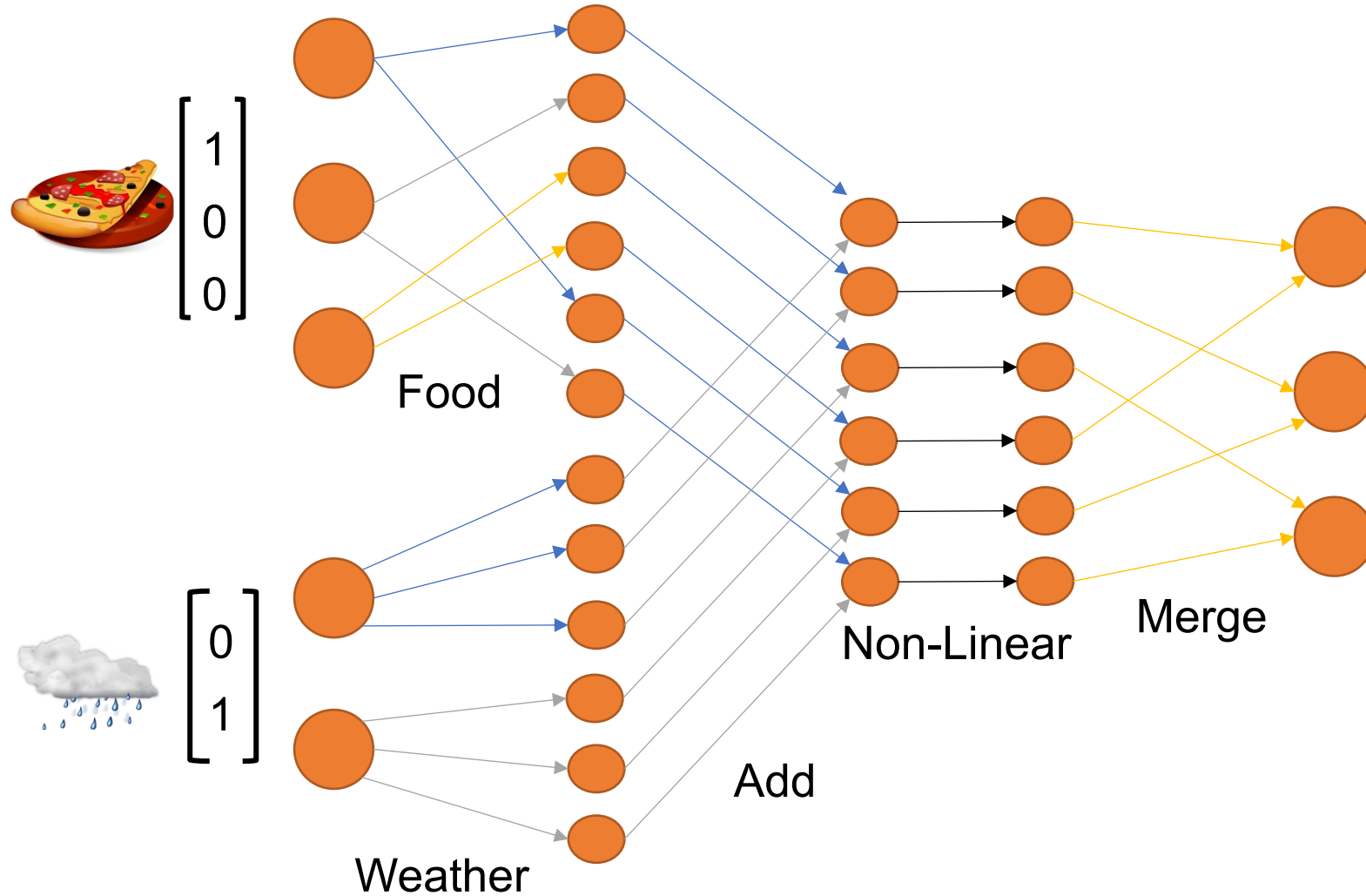
$$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

# Example

## Recurrent neural network
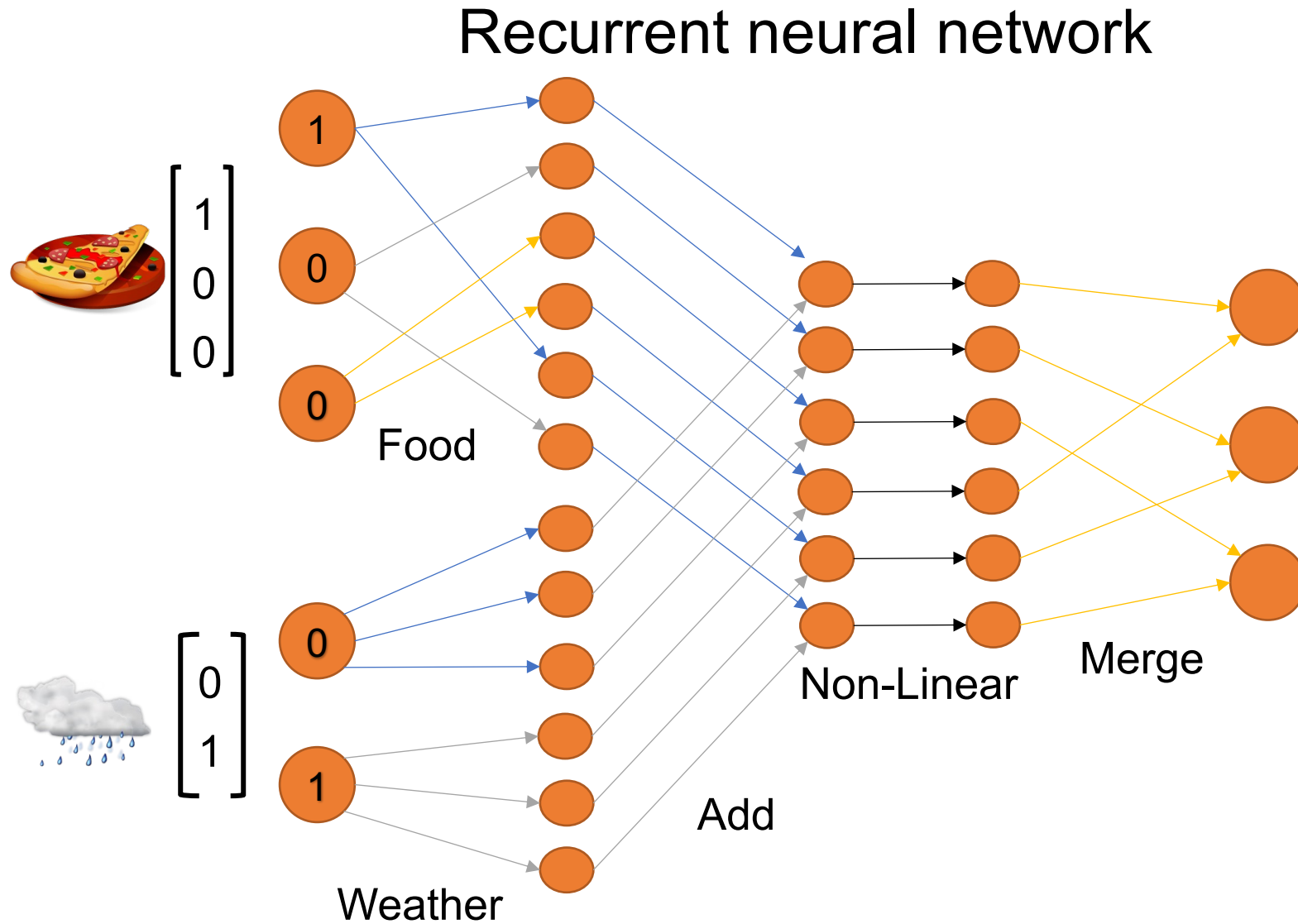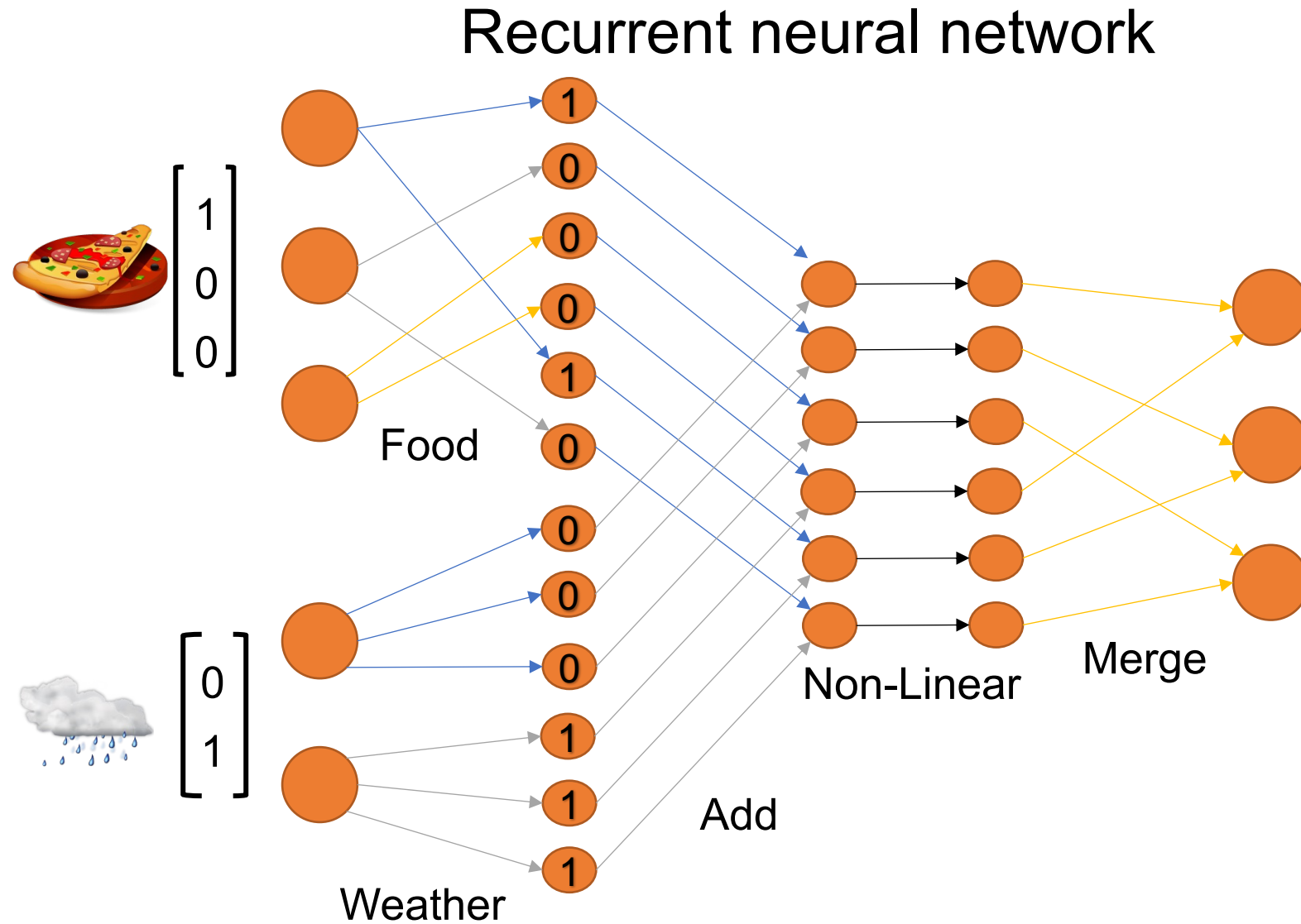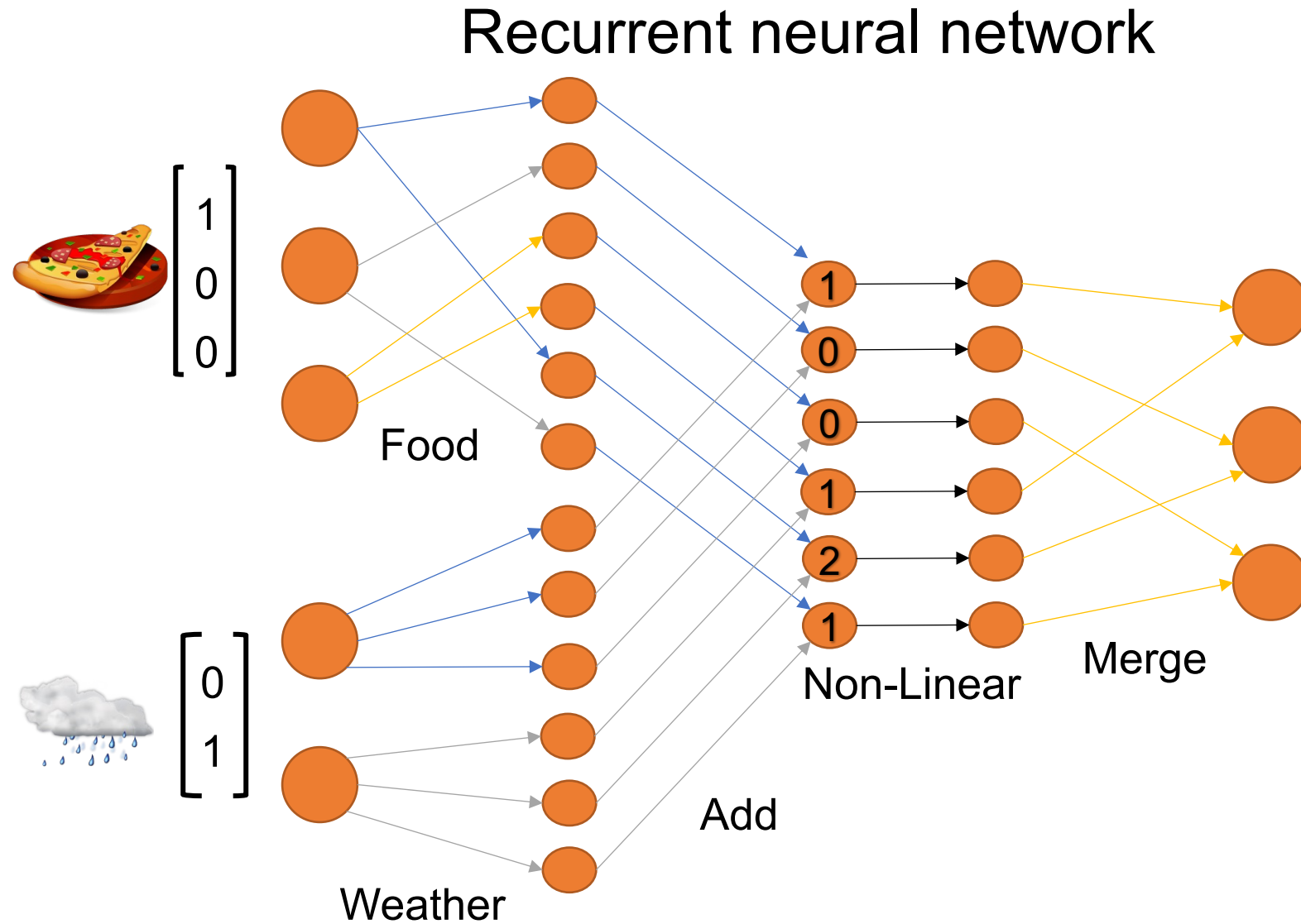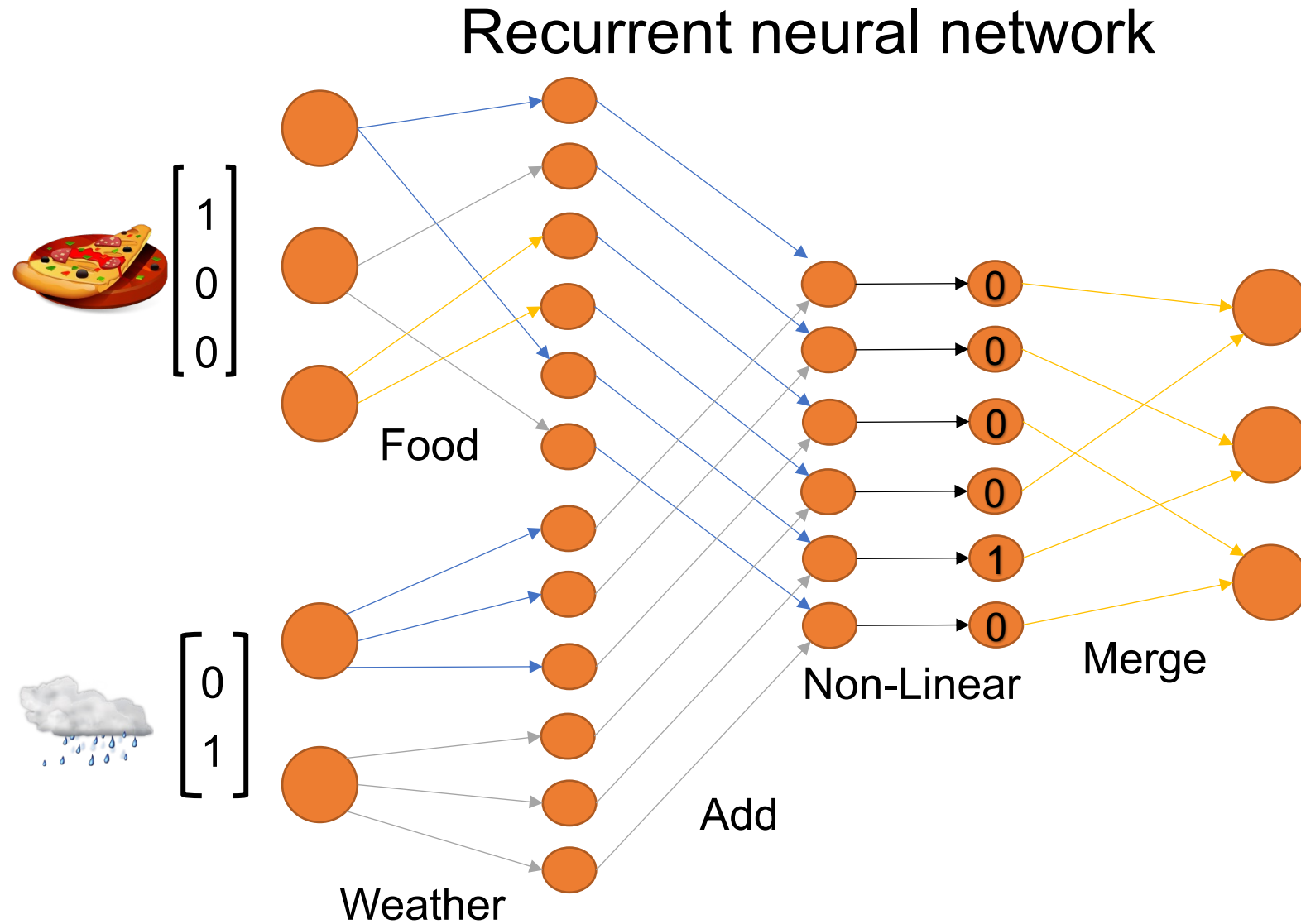


Food

Weather

Add

Non-Linear

Merge

# RNN

## Recurrent neural network



X: inputs

tanh

softmax

H: internal state

Y: outputs

# RNN

## Recurrent neural network



X: inputs

tanh

softmax

H: internal state

Y: outputs

# RNN

_concatenation_

$$X = X_t \mid H_{t-1}$$

$$H_t = \tanh(X.W_H + b_H)$$

$$Y_t = \text{softmax}(H_t.W + b)$$

# Deep RNN

# Long Term Dependency

**Michel C. was born in Paris, France.** He is married and has three children. He received a M.S. in neurosciences from the University Pierre & Marie Curie and the Ecole Normale Supérieure in 1987, and and then spent most of his career in Switzerland, at the Ecole Polytechnique de Lausanne. He specialized in child and adolescent psychiatry and his 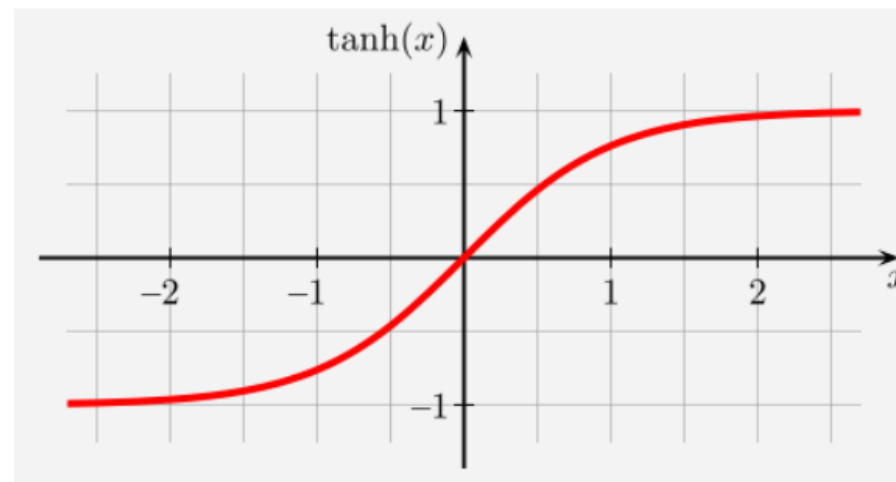first field of research was severe mood disorders in adolescent, topic of his PhD in neurosciences (2002). **His mother tongue is** *? ? ? ? ?*

*Short context* → *English, German, Russian, French …*



*Long context* → *Problems…*

# Simple RNN    Vs. LSTM

# LSTM



- Each line carries an entire vector, from the output of one node to the inputs of others
- The pink circles represent pointwise operations, like vector addition
- The yellow boxes are learned neural network layers.
- Lines merging denote concatenation
- Line forking denote its content being copied and the copies going to different locations

# The Core Idea Behind LSTMs



- The cell state is kind of like a conveyor belt. It runs straight down the entire chain, with only some minor linear interactions.



- Gates are a way to optionally let information through. They are composed out of a sigmoid neural net layer and a pointwise multiplication operation.

- A value of zero means "let nothing through," while a value of one means "let everything through!"

# The Core Idea Behind LSTMs



$$f_t = \sigma \left( W_f \cdot [h_{t-1}, x_t] + b_f \right)$$

- The first step is to decide what information we're going to throw away from the cell state. This decision is made by a sigmoid layer called the "forget gate layer."

- A 1 represents "completely keep this" while a 0 represents "completely get rid of this."

# The Core Idea Behind LSTMs



$$i_t = \sigma\left(W_i \cdot [h_{t-1}, x_t] + b_i\right)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

- It's now time to update the old cell state, $C_{t-1}$, into the new cell state $C_t$. The previous steps already decided what to do, we just need to actually do it.

- We multiply the old state by $f_t$ , forgetting the things we decided to forget earlier. Then we add $i_t * \tilde{C}_t$. This is the new candidate values, scaled by how much we decided to update each state value.



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

# The Core Idea Behind LSTMs



$$o_t = \sigma \left( W_o \left[ h_{t-1}, x_t \right] + b_o \right)$$

$$h_t = o_t * \tanh \left( C_t \right)$$

- Finally, we need to decide what we're going to output. This output will be based on our cell state, but will be a filtered version. First, we run a sigmoid layer which decides what parts of the cell state we're going to output. Then, we put the cell state through tanh and multiply it by the output of the sigmoid gate, so that we only output the parts we decided to.

# LSTM

concatenate : $X = X_t \mid H_{t-1}$

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 2 \\ 1 \end{bmatrix} \xrightarrow{\text{Non-Linear}} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

forget gate : $f = \sigma(X.W_f + b_f)$   $n$

update gate : $u = \sigma(X.W_u + b_u)$   $n$

result gate : $r = \sigma(X.W_r + b_r)$   $n$

input : $X' = \tanh(X.W_c + b_c)$   $n$

new $C$ : $C_t = f * C_{t-1} + u * X'$   $n$

new $H$ : $H_t = r * \tanh(C_t)$   $n$

$X_t$

$H_{t-1}$   $H_t$

$C_{t-1}$   $C_t$

$Y_t$

output : $Y_t = \text{softmax}(H_t.W + b)$   $m$

# Recurrent Neural Networks (RNN)

## Limitations of RNN networks

- **Length of the sentence**

I am from France

I am from France, and I'm a student in computer science at the University of Angers in France.

- **RNNs networks are recurrent unfortunately sentences are processed sequentially word by word**

NN block

↓

I am from France

Time

Calculation and machine power

**Transformer architecture**

While processing a word, Attention enables the model to focus on other words in the input that are closely related to that word.

I am from France and I'm PhD student in computer science in University of Angers in France.
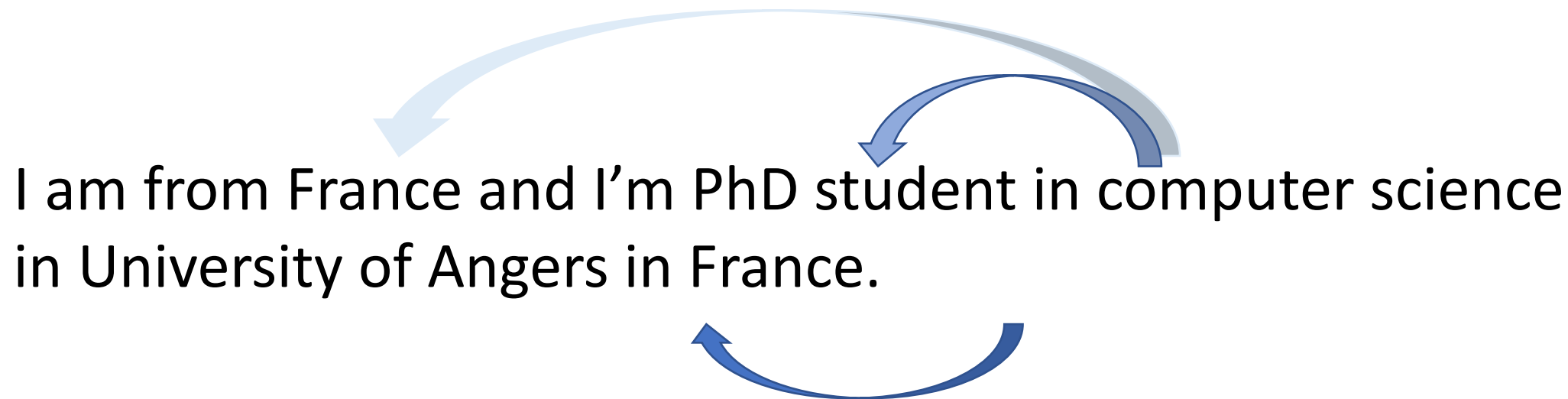
# Architecture

# Architecture

**ENCODER**

Add &Norm

Feed-forward neural network

Add &Norm

$z_1$ $z_2$ $z_3$ $z_4$

Multi-head Self-Attention

$x_1$ $x_2$ $x_3$ $x_4$

**Positional encoding** $+$ $+$ $+$ $+$

$x_1$ $x_2$ $x_3$ $x_4$

I    am    from    France

I    am    from    France

$x_1$ $x_2$ $x_3$ $x_4$ **Embeddings**

$+$ $+$ $+$ $+$

$t_1$ $t_2$ $t_3$ $t_4$ **Positional encoding**

$=$ $=$ $=$ $=$

$x_1$ $x_2$ $x_3$ $x_4$ **Embedding With time signal**

* To give the model a sense of the order of the words, we add positional encoding vectors –

* the values of which follow a specific pattern.

70

# Architecture : Encoder

## Self-attention

**ENCODER**

Add &Norm

Feed-forward neural network

Add &Norm

$z_1$ $z_2$ $z_3$ $z_4$
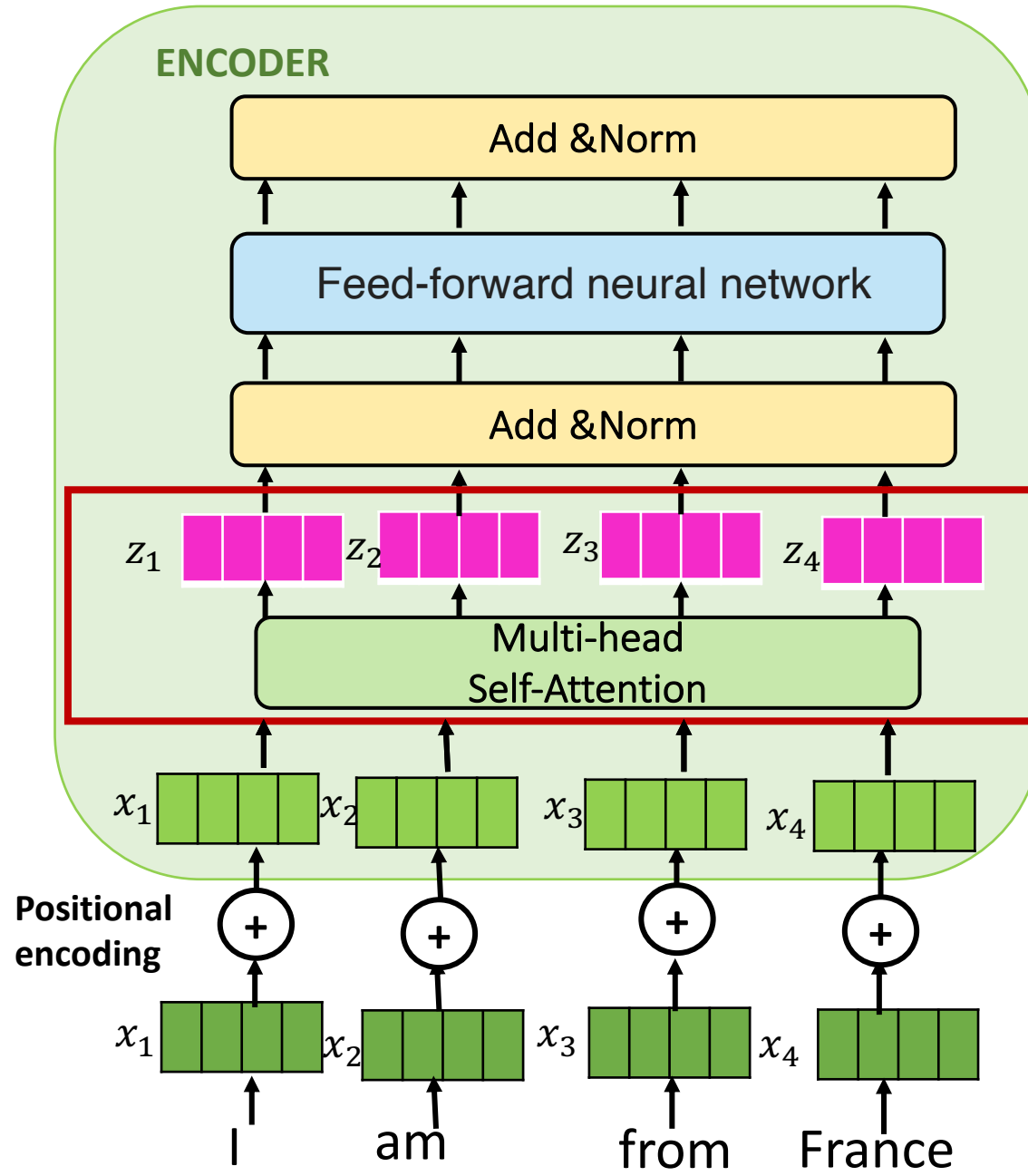
Multi-head Self-Attention

$x_1$ $x_2$ $x_3$ $x_4$

**Positional encoding**

$+$ $+$ $+$ $+$

$x_1$ $x_2$ $x_3$ $x_4$

I    am    from    France

**INPUT Embedding**

I    am

$x_1$    $x_2$

$W^q$

**QUERY**    $q_1$    $q_2$

$W^k$

**KEY**    $k_1$    $k_2$

$W^v$

**VALUE**    $v_1$    $v_2$

**Step 1 :**

- Create a **Query** vector, a **Key** vector, and a **Value** vector.

- These vectors are created by multiplying the embedding by three matrices that we trained during the training process.

71

# Architecture : Encoder



## Self-attention

**Step 2 :.**

- Score each word of the input sentence against this word

- The **score** determines how much focus to place on other parts of the input sentence as we encode a word at a certain position.

72

# Architecture : Encoder

Self-attention

**ENCODER**

Add &Norm

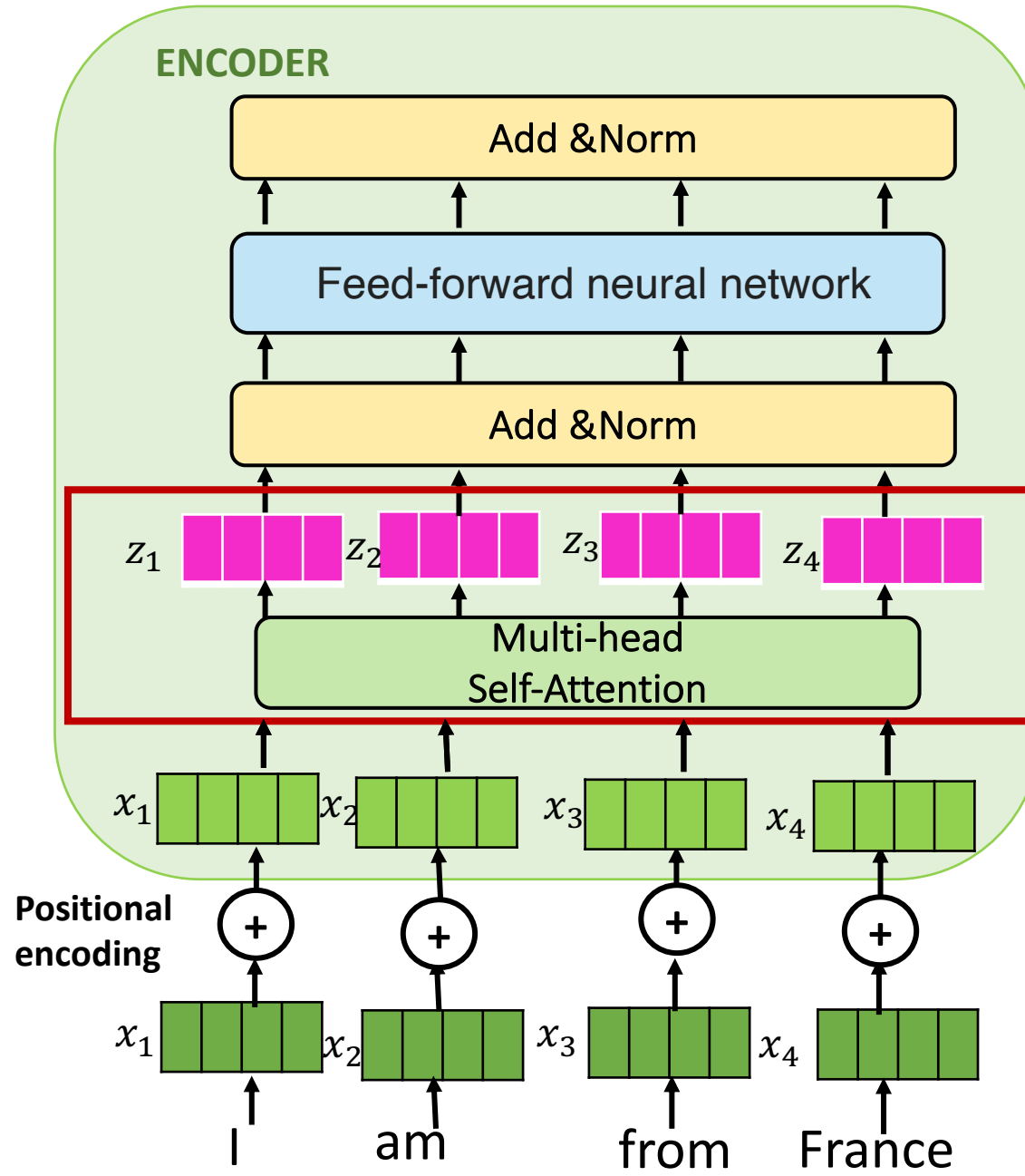Feed-forward neural network

Add &Norm

$z_1$ $z_2$ $z_3$ $z_4$

Multi-head Self-Attention

$x_1$ $x_2$ $x_3$ $x_4$

**Positional encoding**

$+$ $+$ $+$ $+$

$x_1$ $x_2$ $x_3$ $x_4$

I  am  from  France

| | I | am |
|---|---|---|
| **INPUT** | | |
| **Embedding** | $x_1$ | $x_2$ |
| **QUERY** | $q_1$ | $q_2$ |
| **KEY** | $k_1$ | $k_2$ |
| **VALUE** | $v_1$ | $v_2$ |
| **SCORE** | $q_1 * k_1 = 112$ | $q_1 * k_2 = 96$ |
| **DVIDIDE BY 8** $(\sqrt{d_k})$ | 14 | 12 |
| **SOFTMAX** | 0,88 | 0,12 |

With $d_k$ : dimension of $x_n$ vector (64 in the original paper)

**Step 3 & 4:**
- **Divide** the scores by 8 then pass the result through a softmax operation.

- **Softmax** score determines how much each word will be expressed at this position.

73

# Architecture : Encoder

## Self-attention

**INPUT**

**Embedding**

**QUERY**
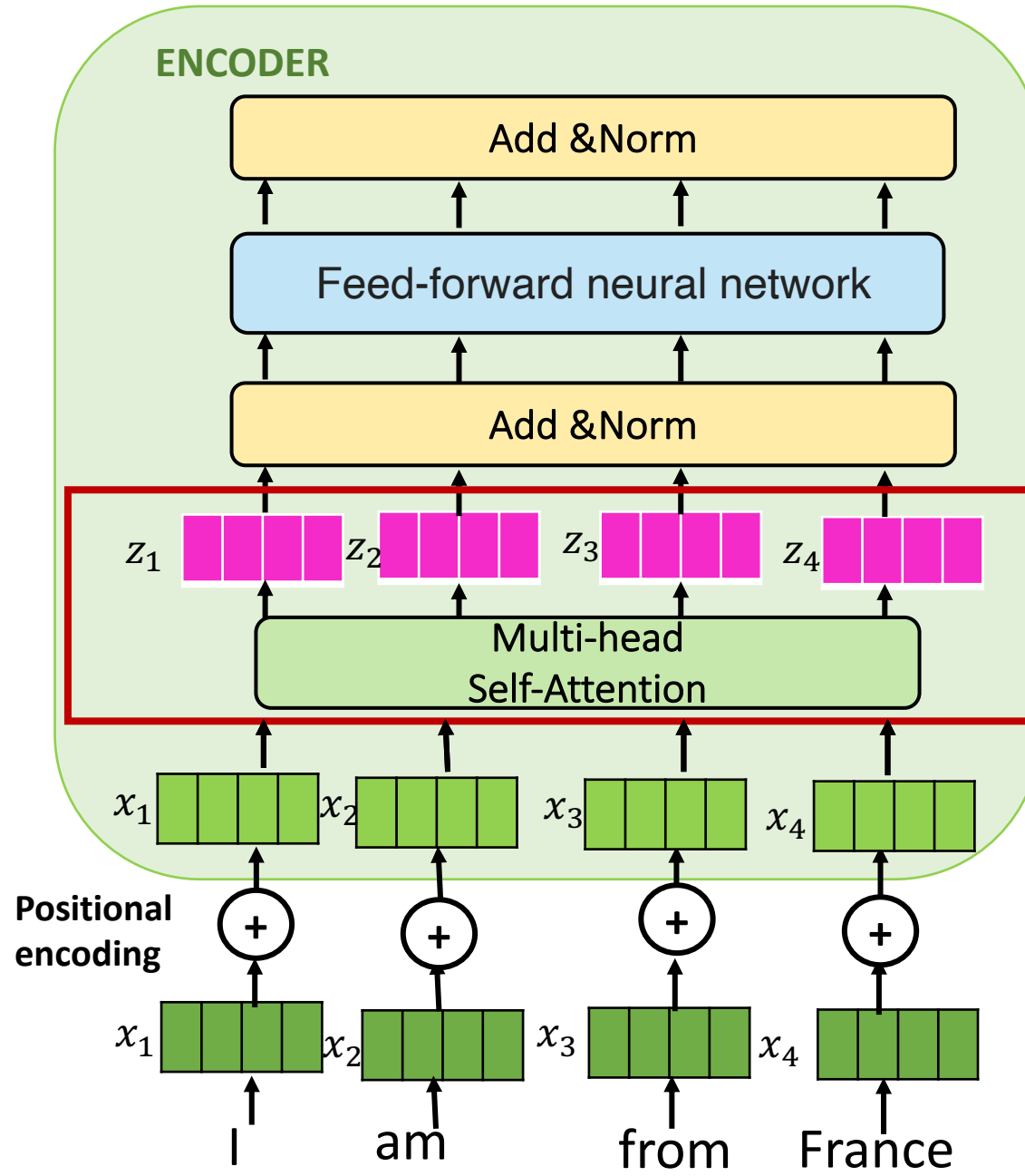
**KEY**

**VALUE**

| | I | am |
|---|---|---|
| Embedding | $x_1$ | $x_2$ |
| QUERY | $q_1$ | $q_2$ |
| KEY | $k_1$ | $k_2$ |
| VALUE | $v_1$ | $v_2$ |
| SCORE | $q_1 * k_1 = 112$ | $q_1 * k_2 = 96$ |
| DVIDIDE BY 8 | 14 | 12 |
| SOFTMAX | 0,88 | 0,12 |
| SOFTMAX * VALUE | $v_1$ | $v_2$ |

**Step 5 :**
- **Multiply** each value vector by the softmax score

# Architecture : Encoder

## Self-attention

**ENCODER**

Add &Norm

Feed-forward neural network

Add &Norm

$z_1$ $z_2$ $z_3$ $z_4$

Multi-head Self-Attention
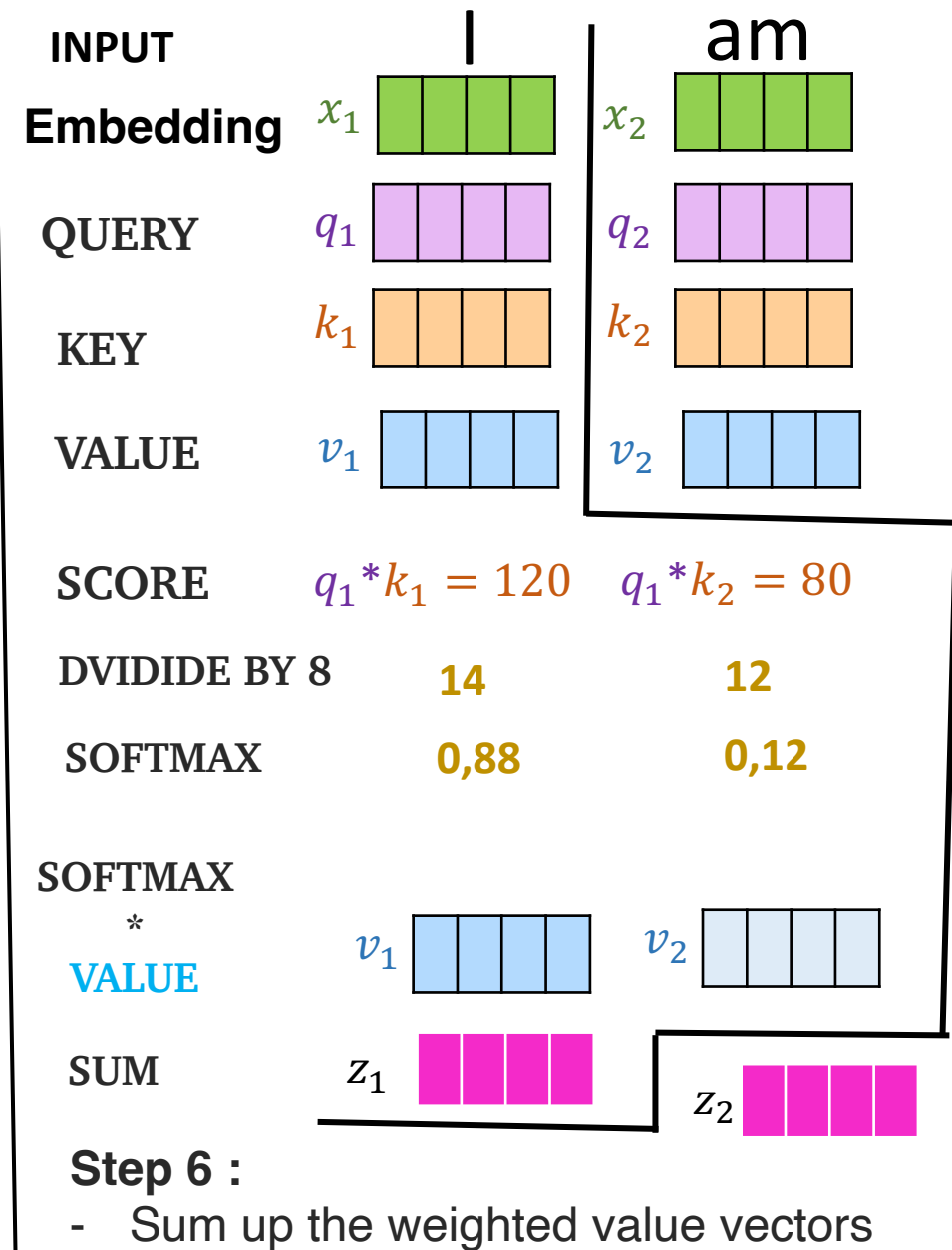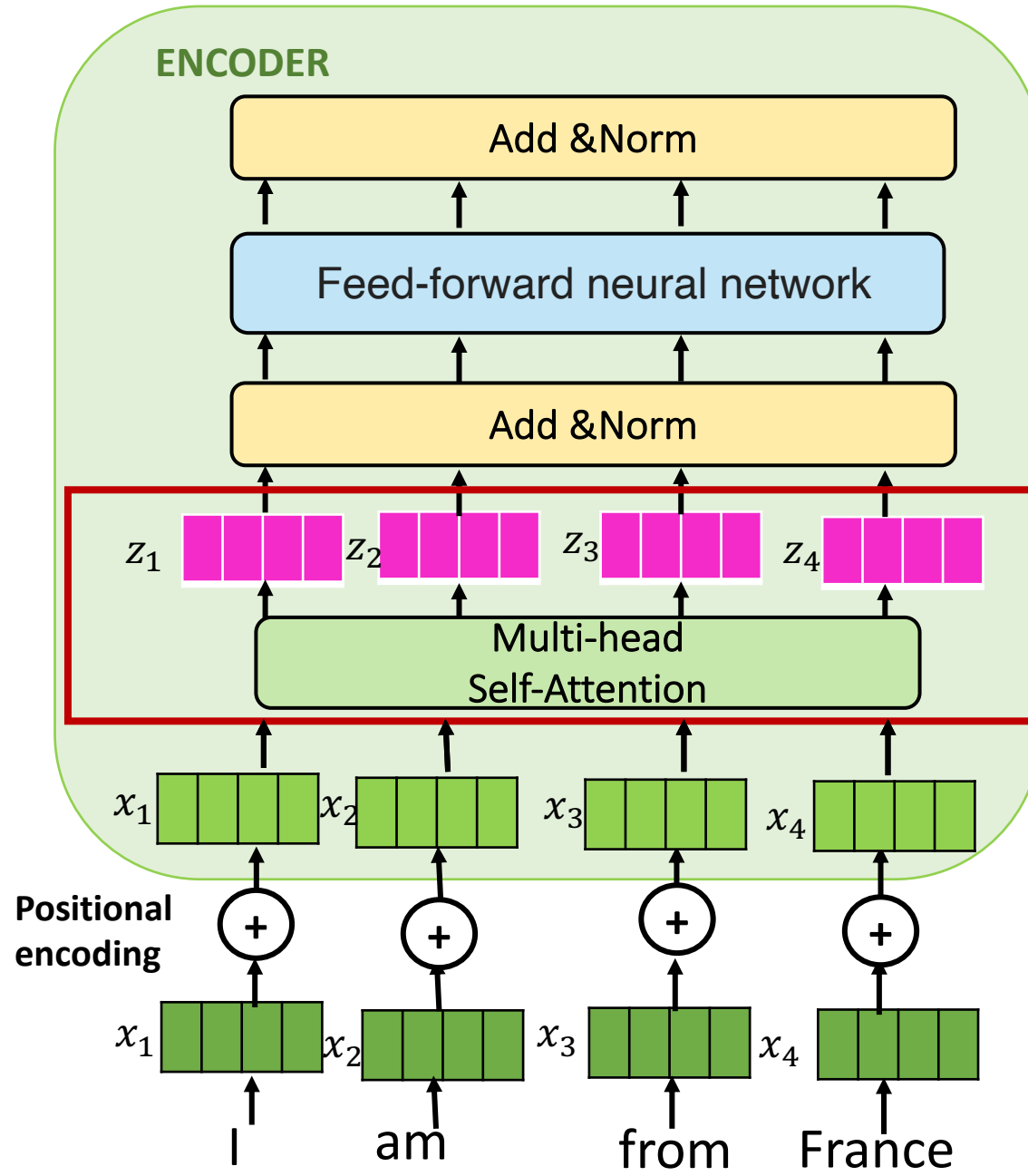
$x_1$ $x_2$ $x_3$ $x_4$

**Positional encoding**

$x_1$ $x_2$ $x_3$ $x_4$

I      am      from      France

| | I | am |
|---|---|---|
| **INPUT** | | |
| **Embedding** | $x_1$ | $x_2$ |
| **QUERY** | $q_1$ | $q_2$ |
| **KEY** | $k_1$ | $k_2$ |
| **VALUE** | $v_1$ | $v_2$ |
| **SCORE** | $q_1 * k_1 = 120$ | $q_1 * k_2 = 80$ |
| **DVIDIDE BY 8** | 14 | 12 |
| **SOFTMAX** | 0,88 | 0,12 |

**SOFTMAX * VALUE**
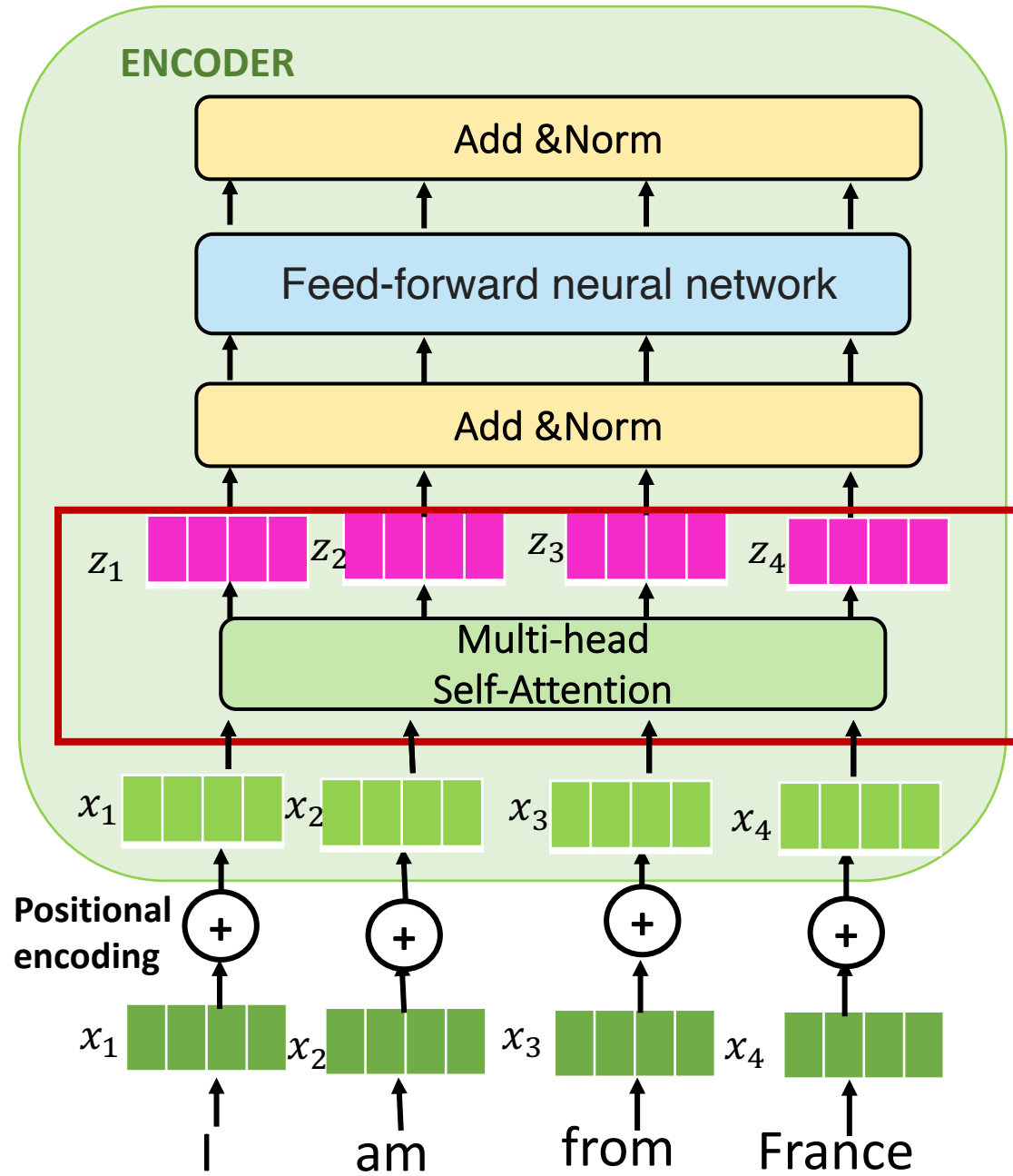
$v_1$ $v_2$

**SUM** $z_1$ $z_2$

**Step 6 :**
- Sum up the weighted value vectors
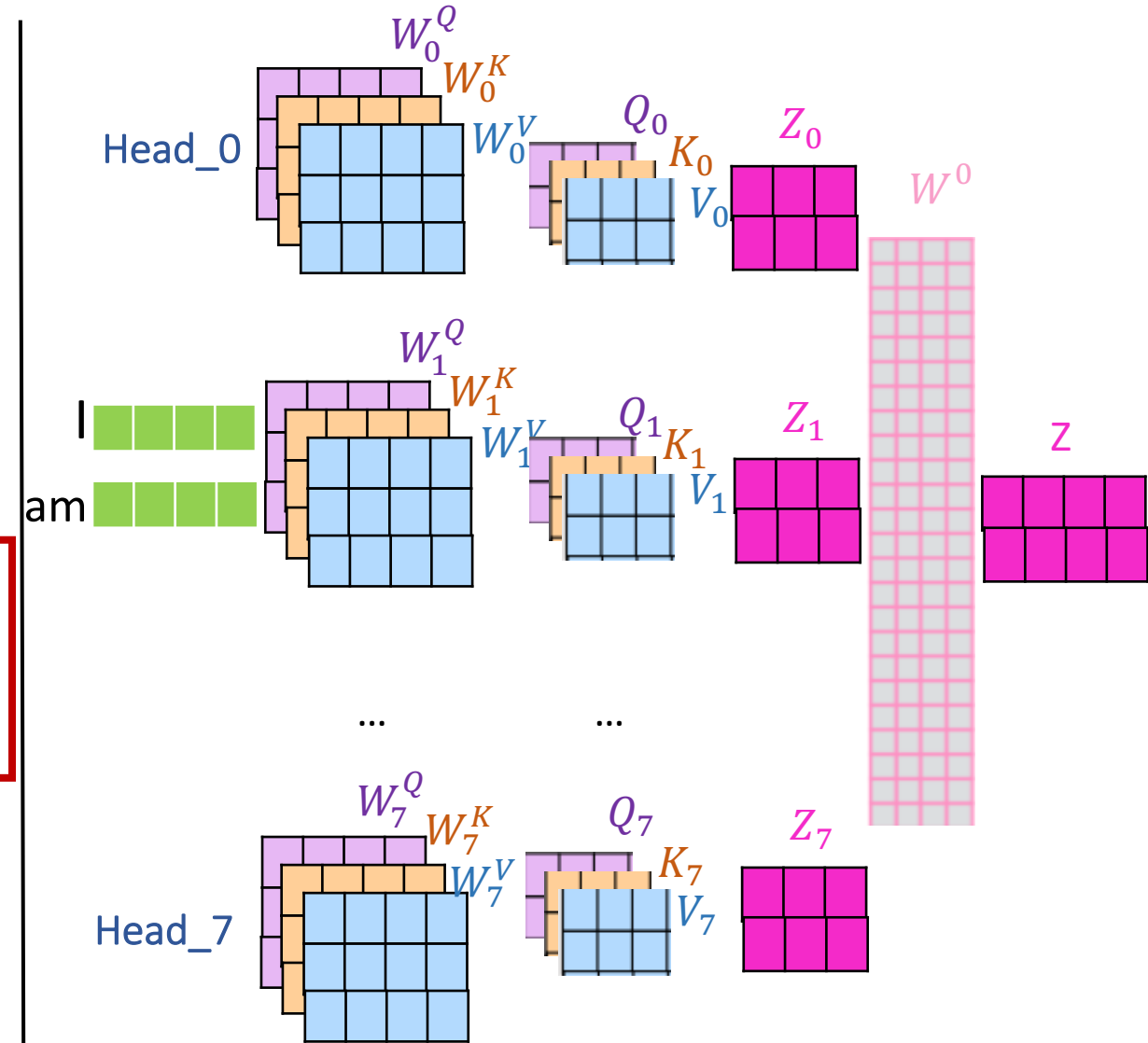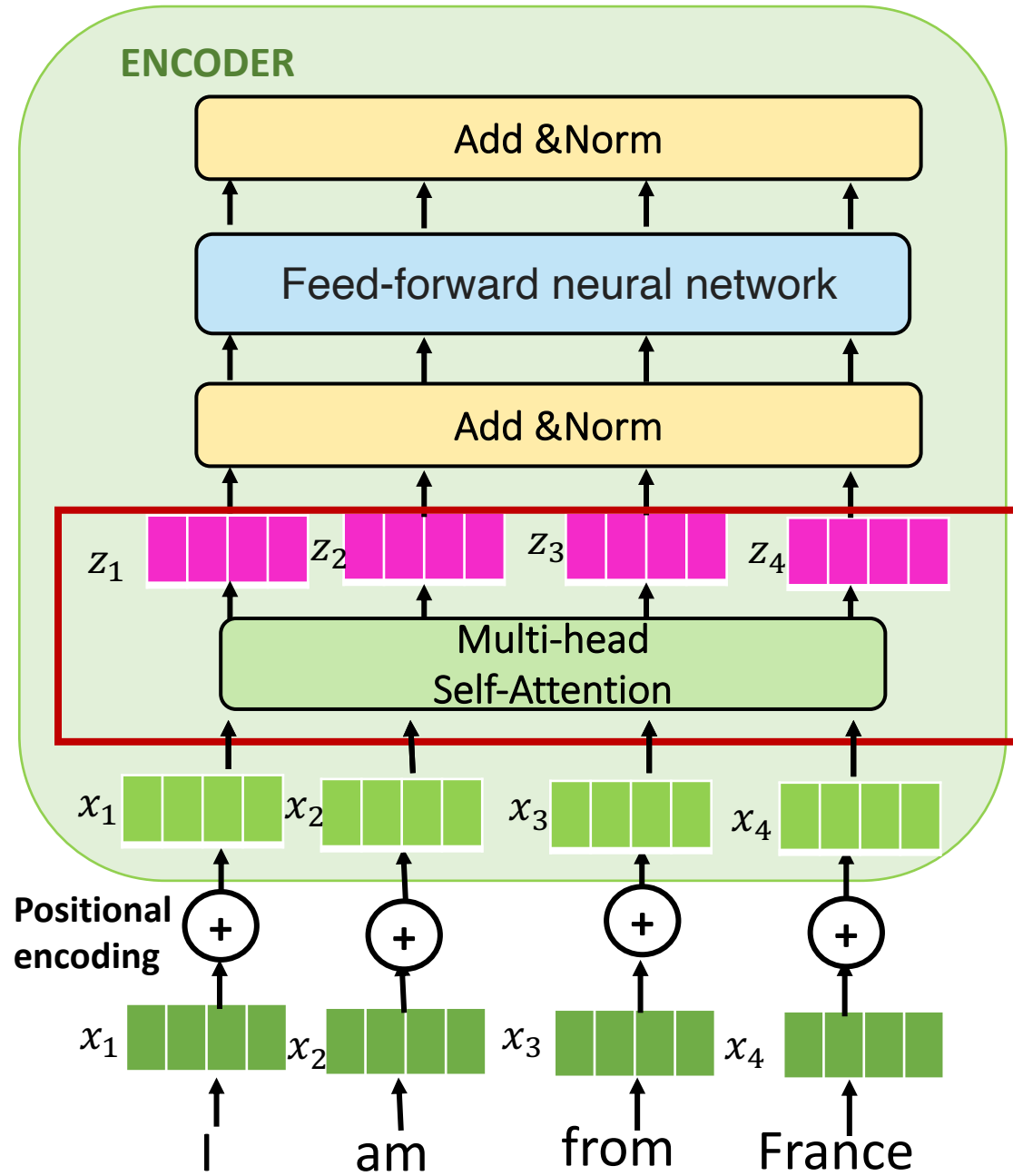
75

# Architecture : Encoder

## Multi-head of self-attention



- Expands the model's ability to focus on different positions

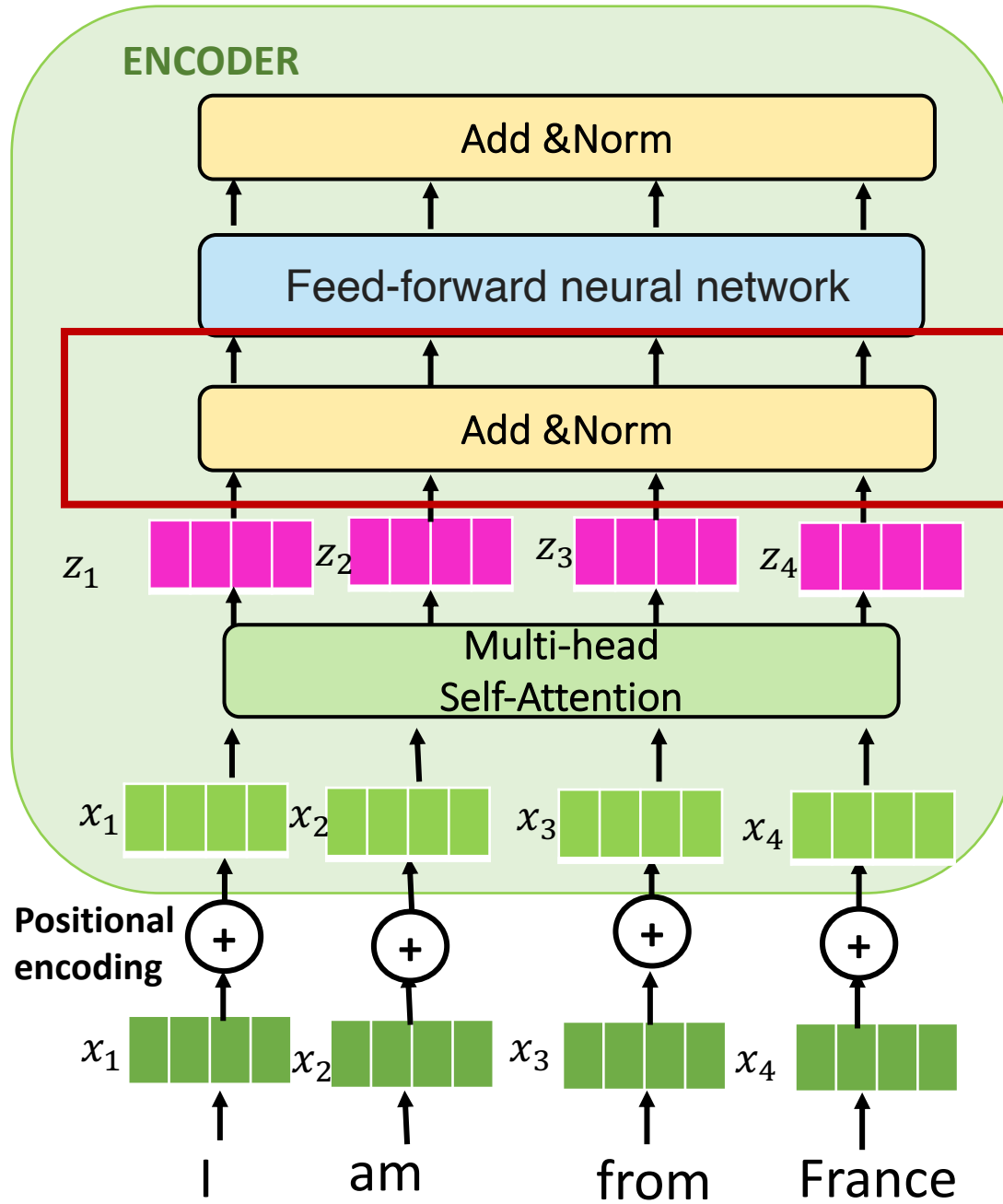- Gives the attention layer multiple "representation subspaces".

# Architecture : Encoder

## Add & Normalisation



$$\text{LayerNorm}(\underset{X}{\boxed{\phantom{aaa}}} + \underset{Z}{\boxed{\phantom{aaa}}})$$

# Architecture : Encoder

## Feed-forward neural network

- **FFNs** are fully connected layers

- **FFN** is a position-wise network

- **FFN** contains two layers and applies a ReLu activation function

**FFN(x) = max(0, xW1 + b1)W2 + b2**

# Architecture : DECODER

## Multi-head of self-attention



- In the decoder, the self-attention layer is only allowed to attend to earlier positions in the output sequence.

- This is done by masking future positions (setting them to -inf) before the softmax step in the self-attention calculation.

## Encoder-Decoder-Attention

| Softmax |
| --- |

| Linear |
| --- |

**DECODER#2**

**DECODER#1**

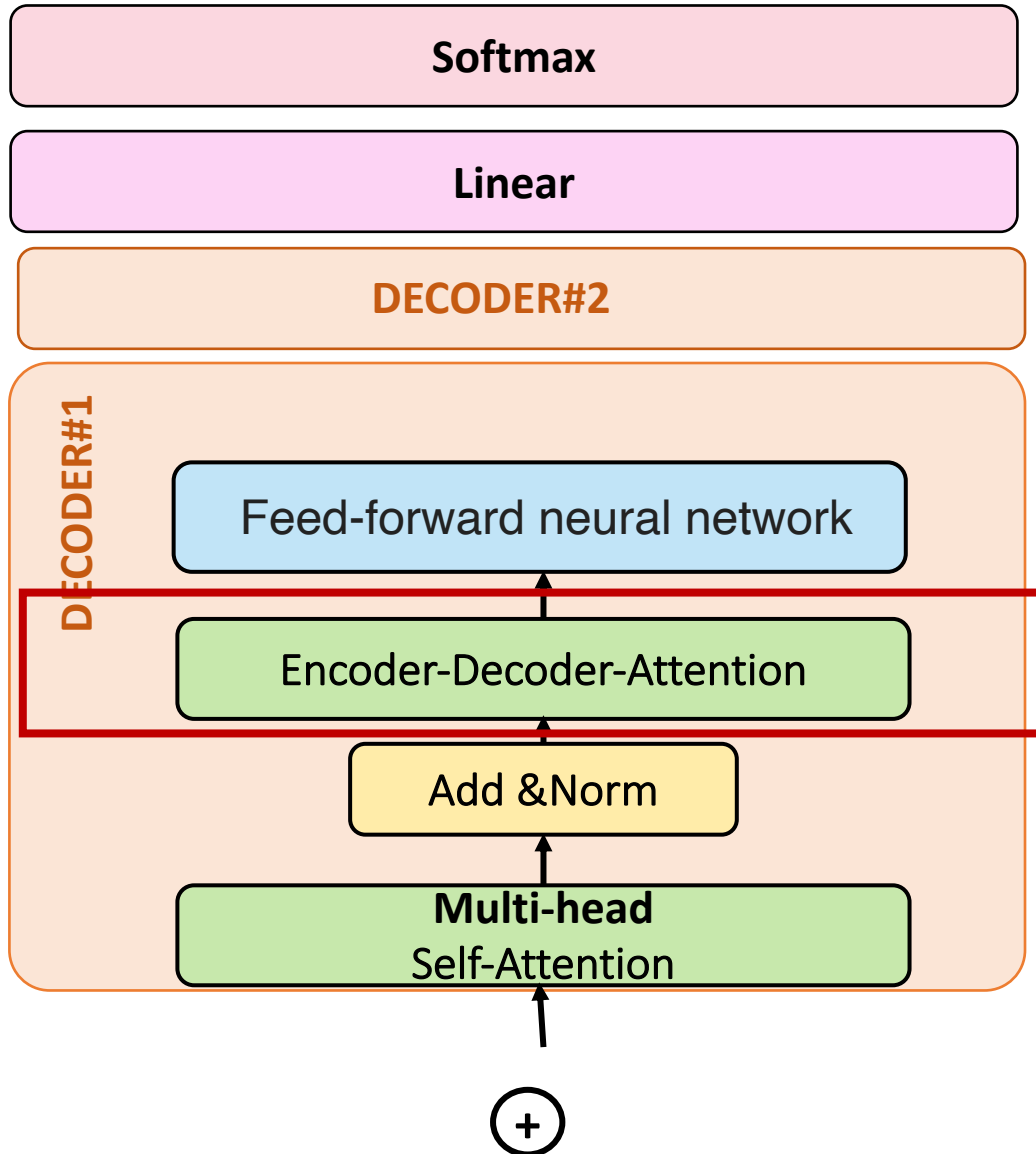| Feed-forward neural network |
| --- |

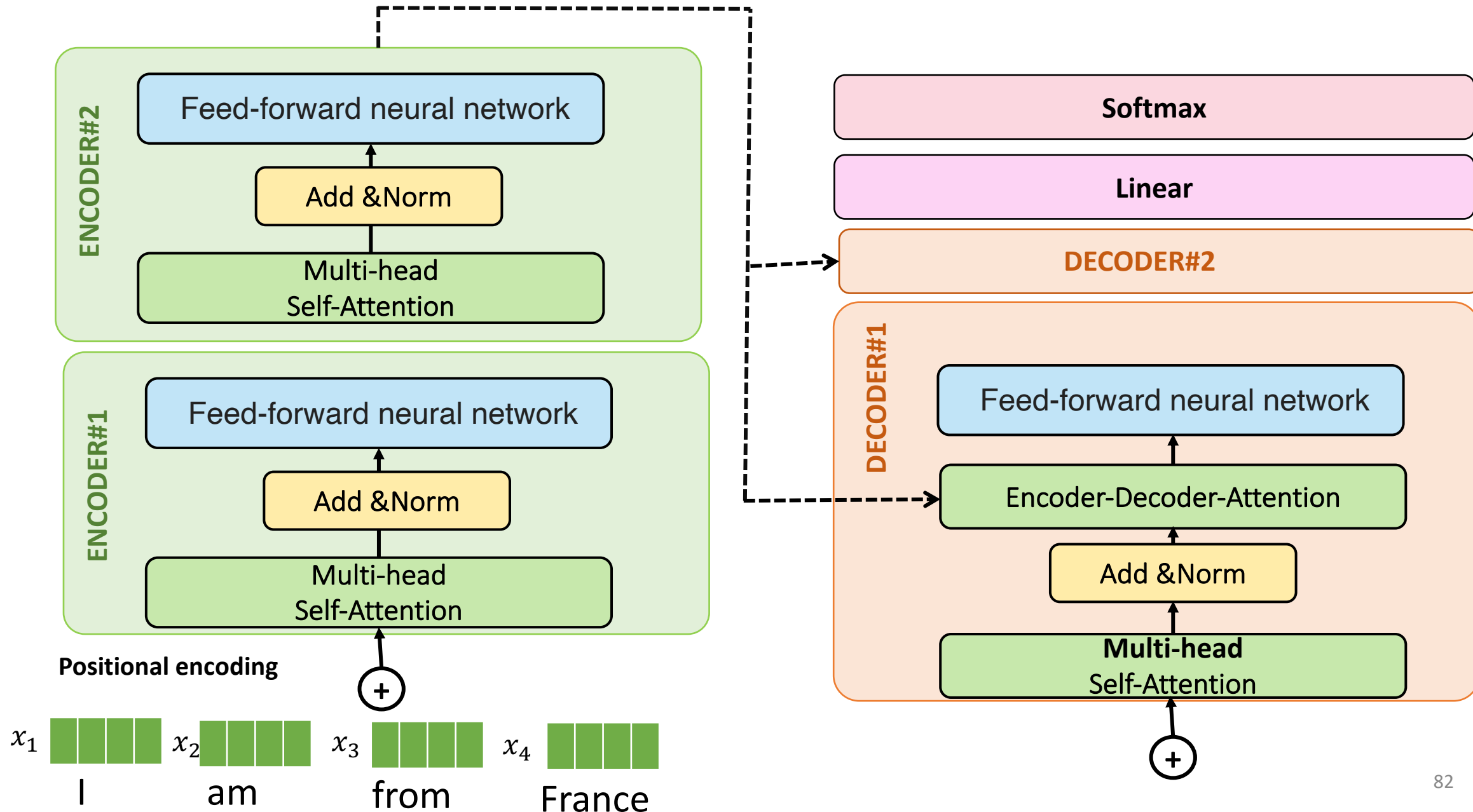| Encoder-Decoder-Attention |
| --- |

| Add &Norm |
| --- |

| **Multi-head** <br> Self-Attention |
| --- |

( + )

- The "Encoder-Decoder Attention" layer works just like multiheaded self-attention, except it creates its Queries matrix from the layer below it, and takes the Keys and Values matrix from the output of the encoder stack.
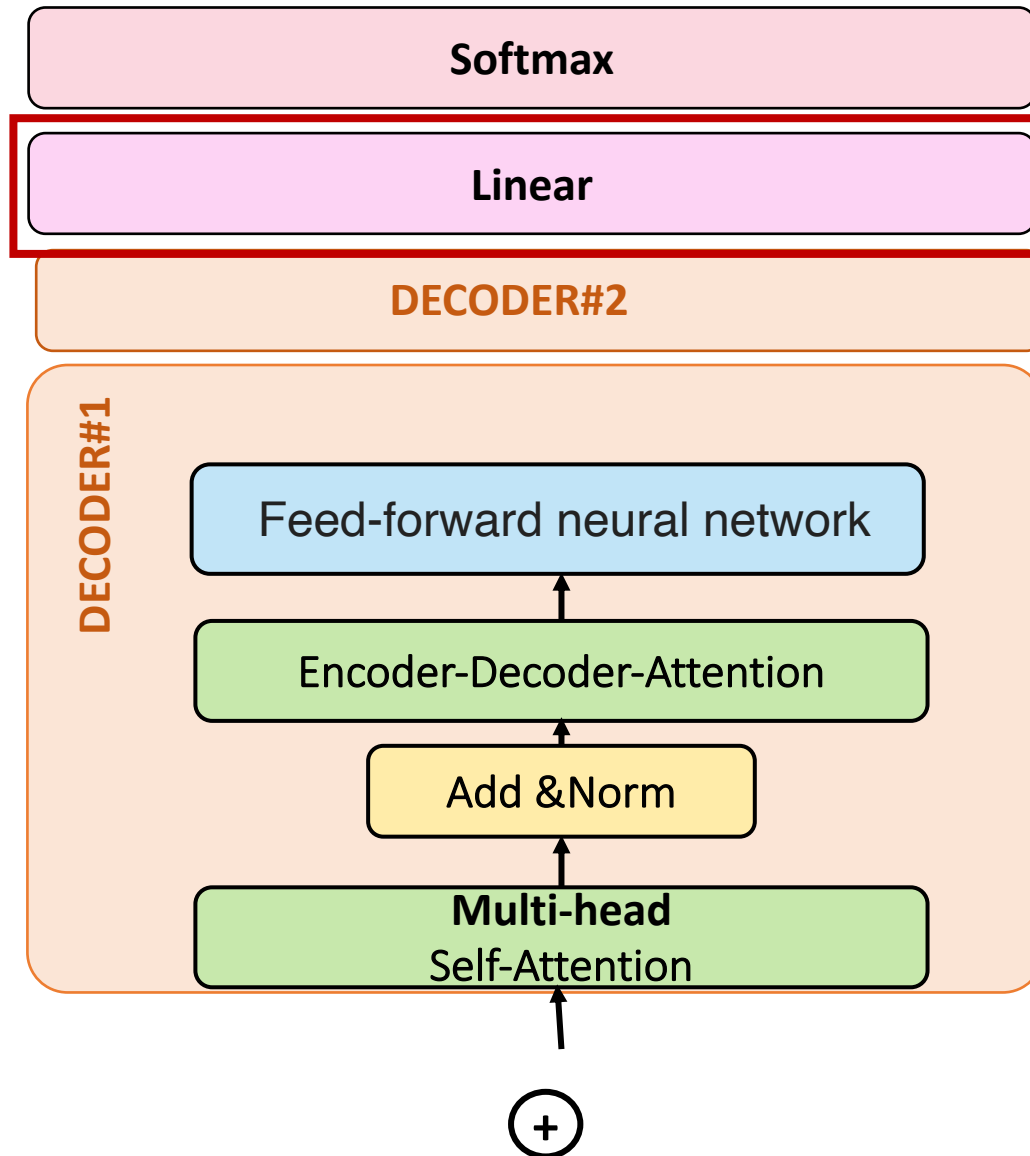
K          V

81

# Architecture

# Architecture : DECODER

## Linear

| Softmax |
|---|

| **Linear** |
|---|

**DECODER#2**

**DECODER#1**

| Feed-forward neural network |
|---|

| Encoder-Decoder-Attention |
|---|

| Add &Norm |
|---|

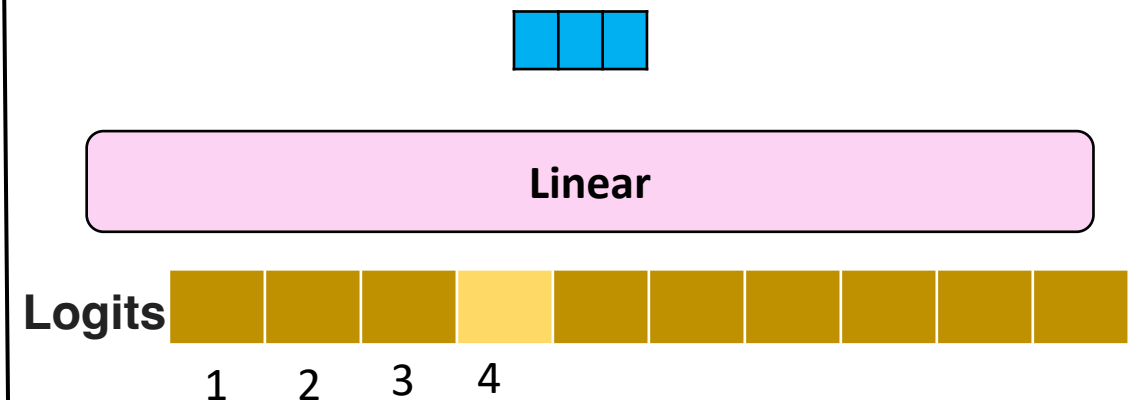| **Multi-head** Self-Attention |
|---|

(+)

**The decoder stack outputs a vector of floats. How do we turn that into a word?**
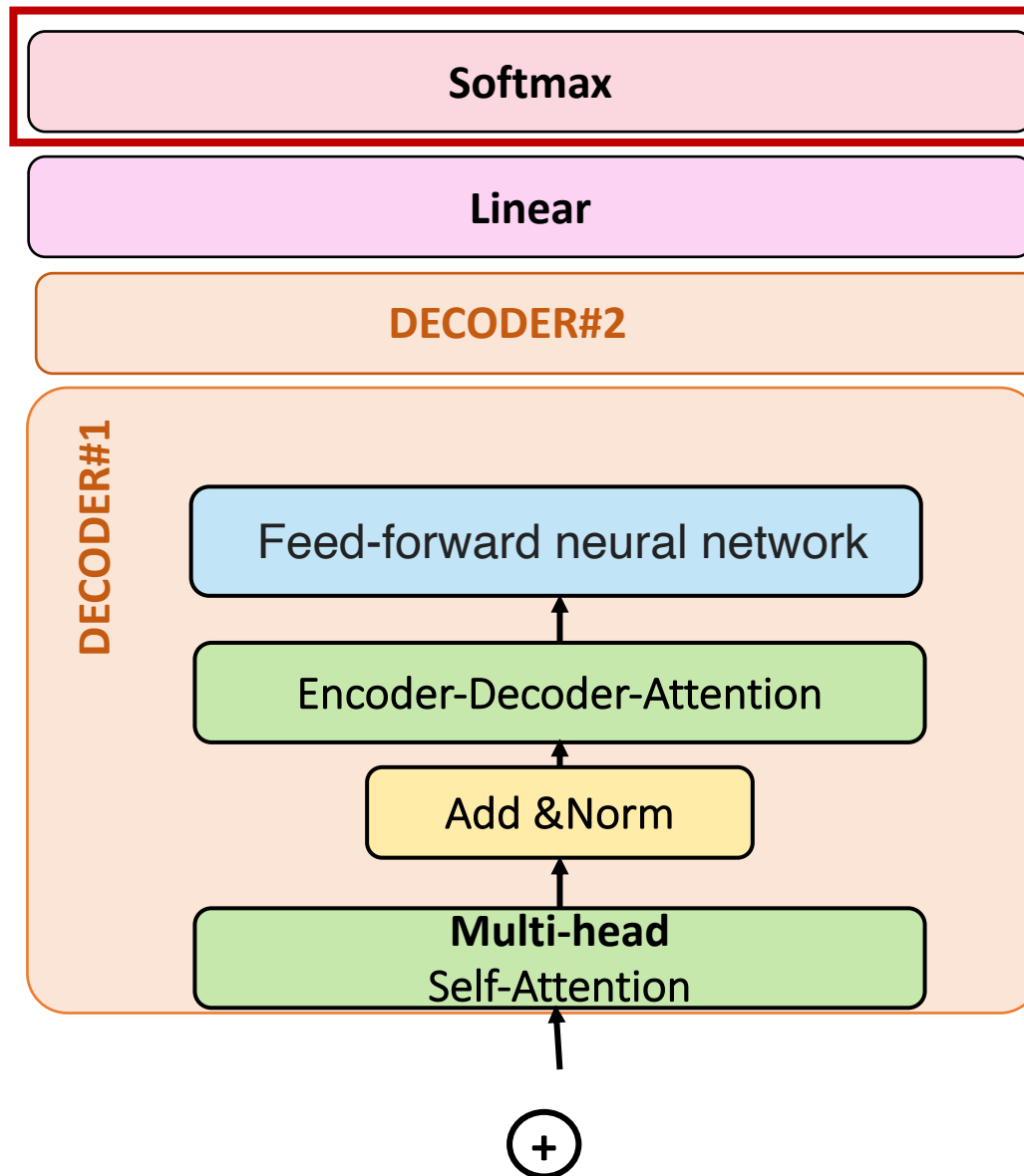
**Word**

- The Linear layer is a simple fully connected neural network that projects the vector produced by the stack of decoders, into a much, much larger vector called a **logits vector.**
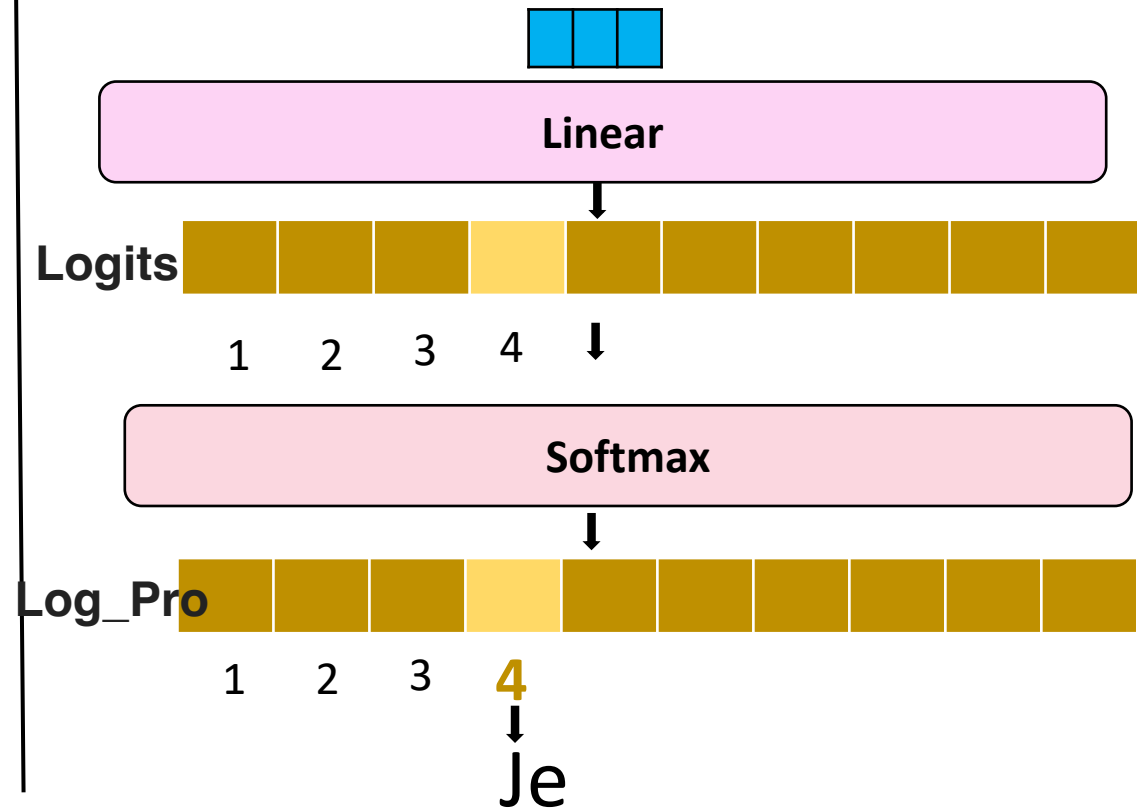
| **Linear** |
|---|

**Logits**

1    2    3    4

# Architecture : DECODER

**Softmax**

**Linear**

**DECODER#2**

**DECODER#1**

Feed-forward neural network

Encoder-Decoder-Attention

Add &Norm

**Multi-head**
Self-Attention
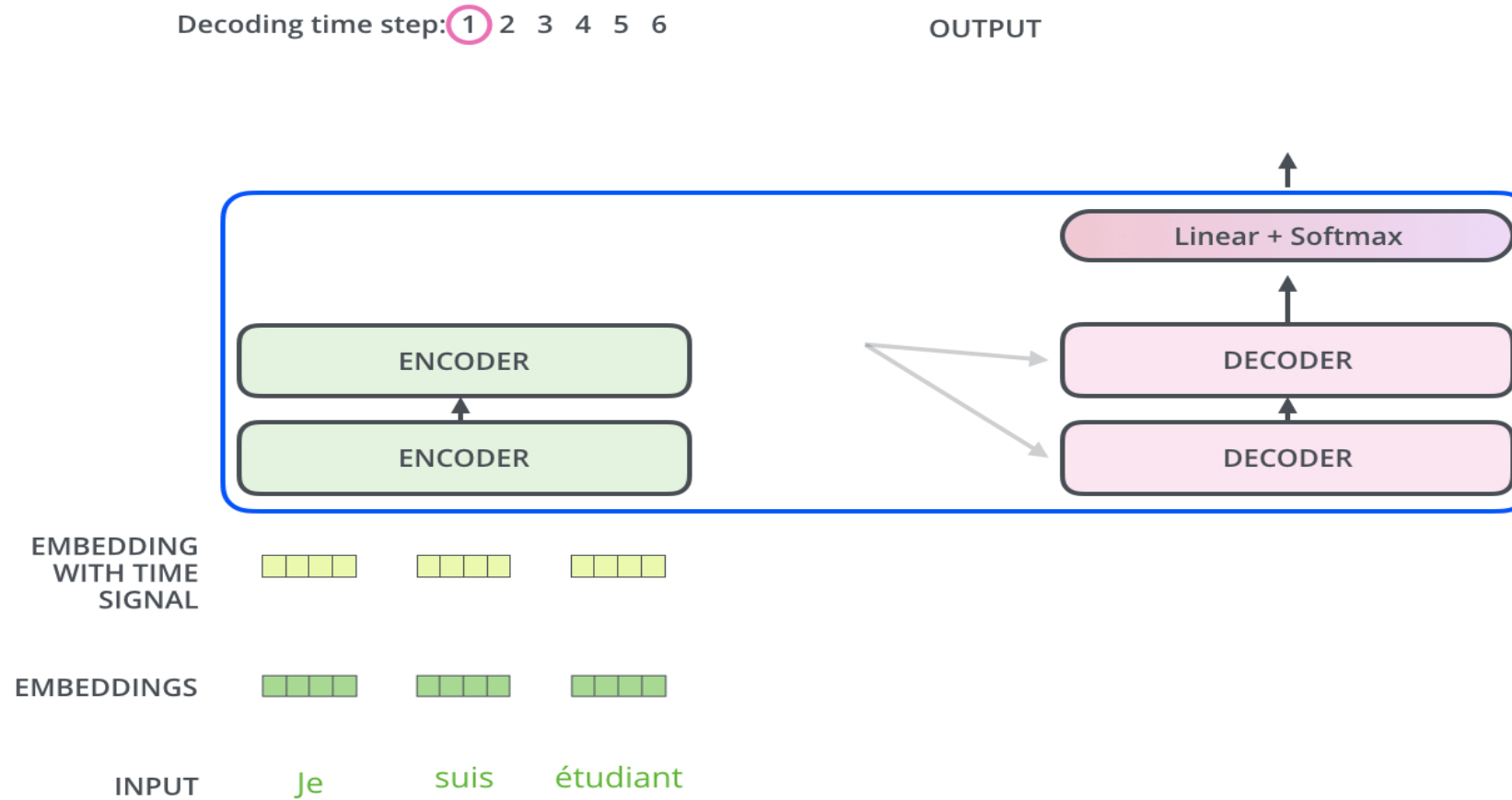
⊕

- The softmax layer then turns those scores into probabilities.

- The cell with the highest probability is chosen, and the word associated with it is produced as the output for this time step.

**Linear**

**Logits**

1   2   3   4

**Softmax**

**Log_Pro**

1   2   3   **4**

Je

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

OUTPUT



EMBEDDING WITH TIME SIGNAL

EMBEDDINGS

INPUT        Je        suis        étudiant
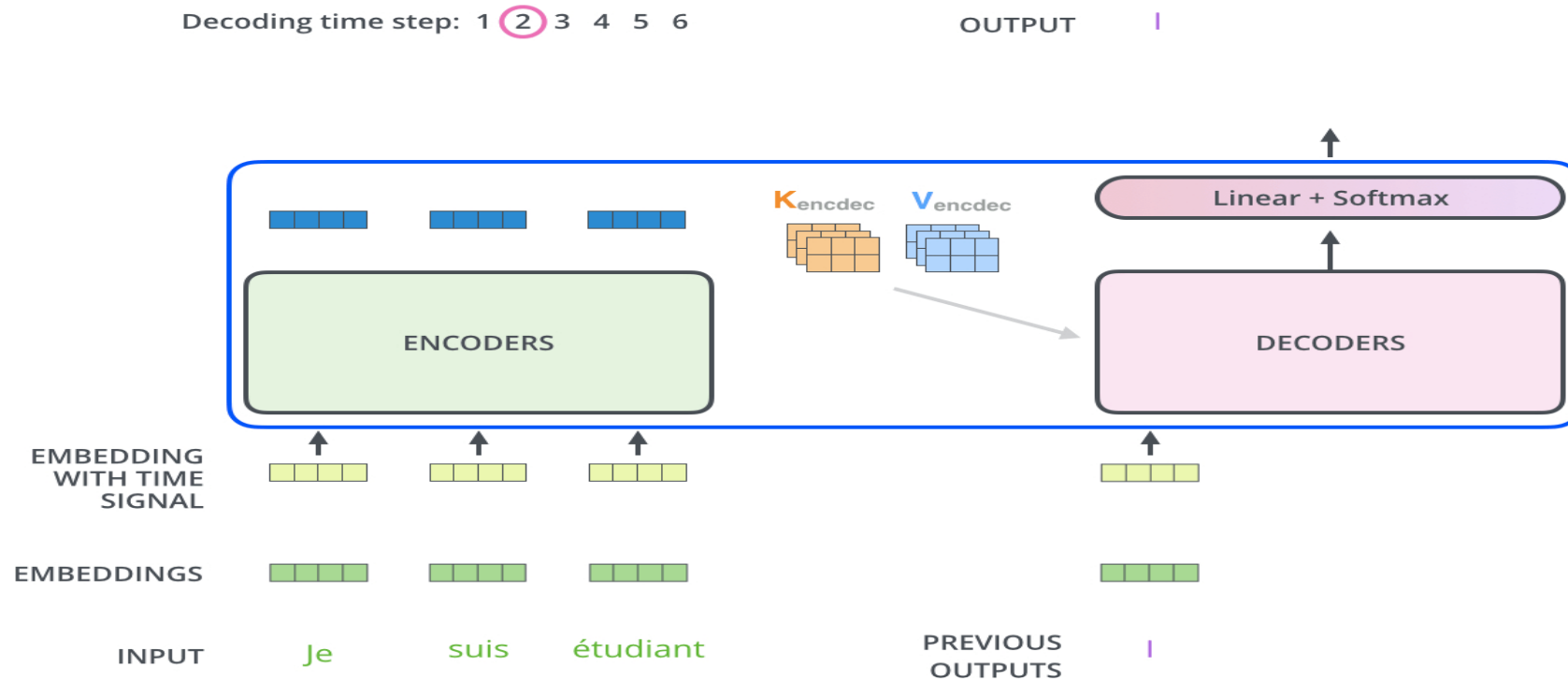
- The encoder start by processing the input sequence.

- The output of the top encoder is then transformed into a set of attention vectors K and V.

# Architecture : ENCODER &DECODER

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

- The output of each step is fed to the bottom decoder in the next time step, and the decoders bubble up their decoding results

- Embed and add positional encoding to those decoder inputs to indicate the position of each word