# Foundations of NLP

CS3126
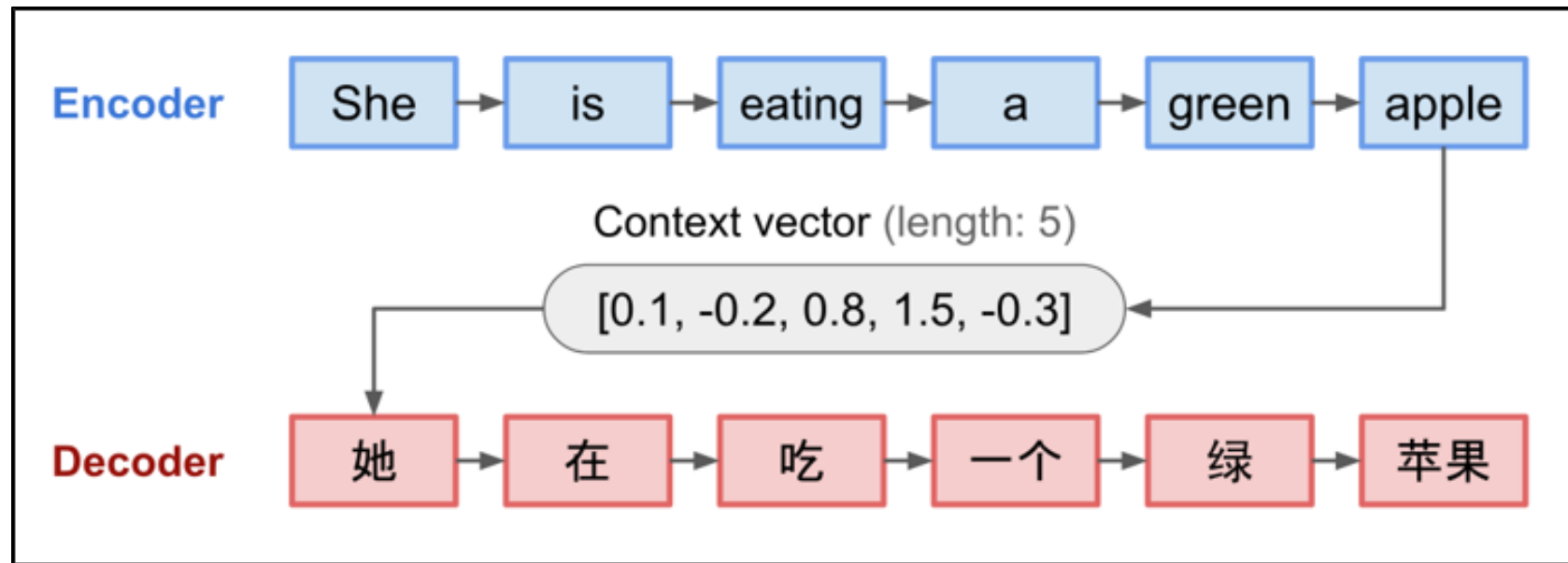Lecture-7
Encoder-decoder models, Attention mechanism

**Mahindra™**
**University** 1
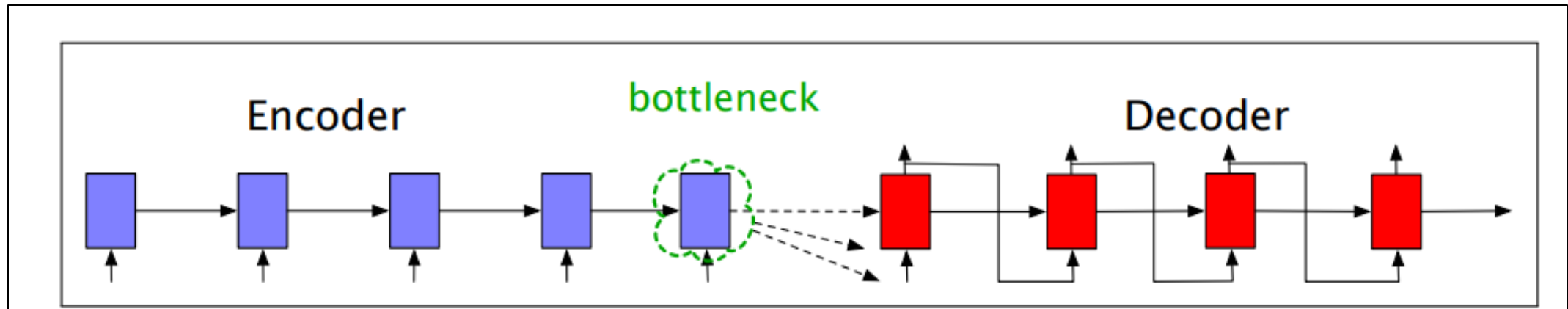Global Thinkers. Engaged Leaders.

# Recap

- Language modeling

- Recurrent Neural Network and Implementation

- Applications of Recurrent Neural Network

- Language modeling using Long Short-term Memory

# Encoder-Decoder model

# Problem- Bottleneck in Encoder-decoder



Requiring the **context c** to be only the encoder's final hidden state forces all the information from the entire source sentence to pass through this representational bottleneck

# Problems with Sequence to Sequence models

- fixed-length context vector design
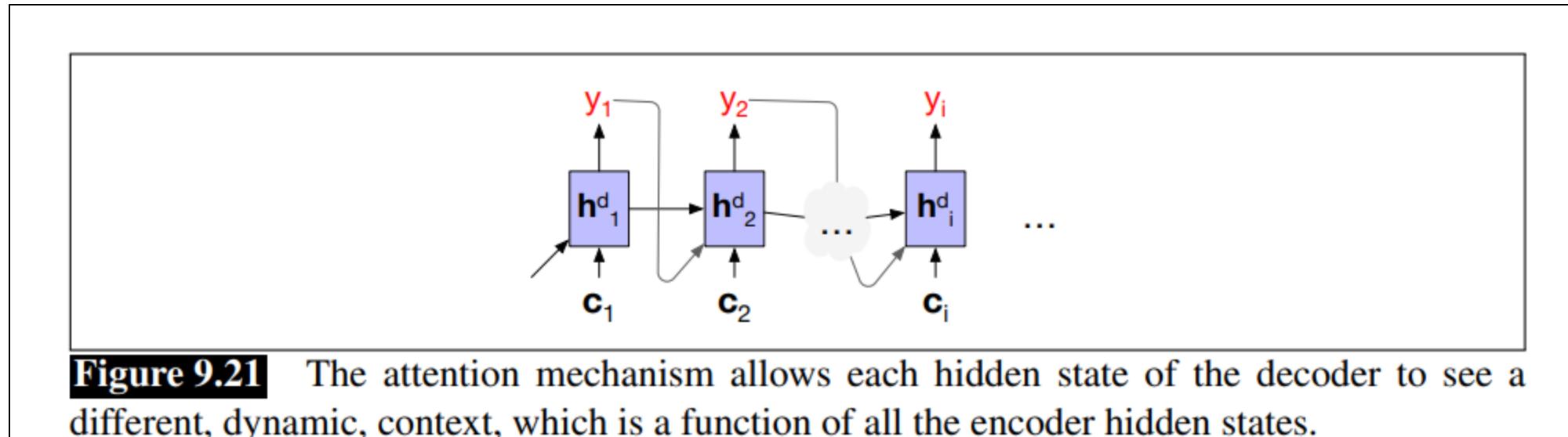  - incapability of remembering long sentences

**Imagine the whole universe in all its beauty - try to visualize everything you can find there and how you can describe it in words. Then imagine all of it is compressed into a single vector of size e.g. 512. Do you feel that the universe is still ok?**

Not only it is hard for the encoder to put all information into a single vector - this is also hard for the decoder.

The decoder sees only one representation of source. However, at each generation step, different parts of source can be more useful than others.

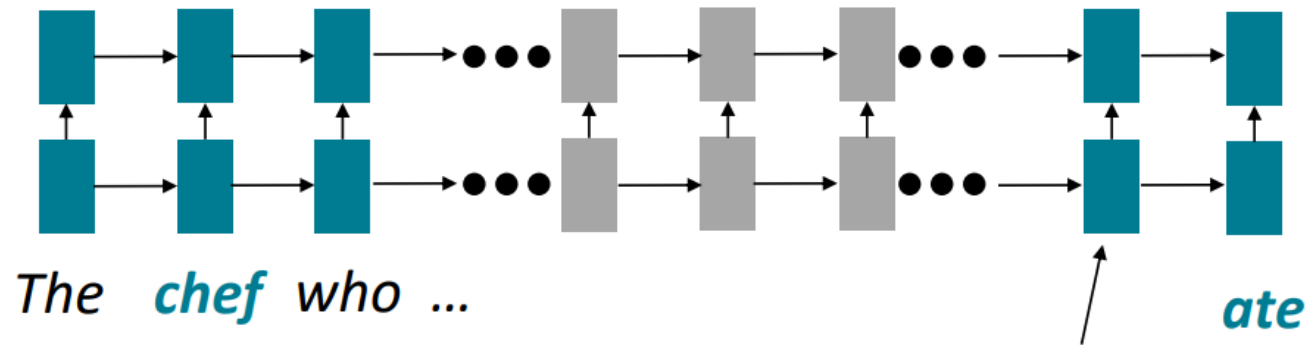https://lena-voita.github.io/nlp_course/seq2seq_and_attention.html

# Solution to bottleneck problem: Attention

Allow the decoder to get information from all the hidden states of the encoder, not just the last hidden state.



**Figure 9.21** The attention mechanism allows each hidden state of the decoder to see a different, dynamic, context, which is a function of all the encoder hidden states.

Image Reference: Speech and Language Processing by Daniel Jurafsky and James H. Martin

# Issues with Recurrent models: Linear interaction distance

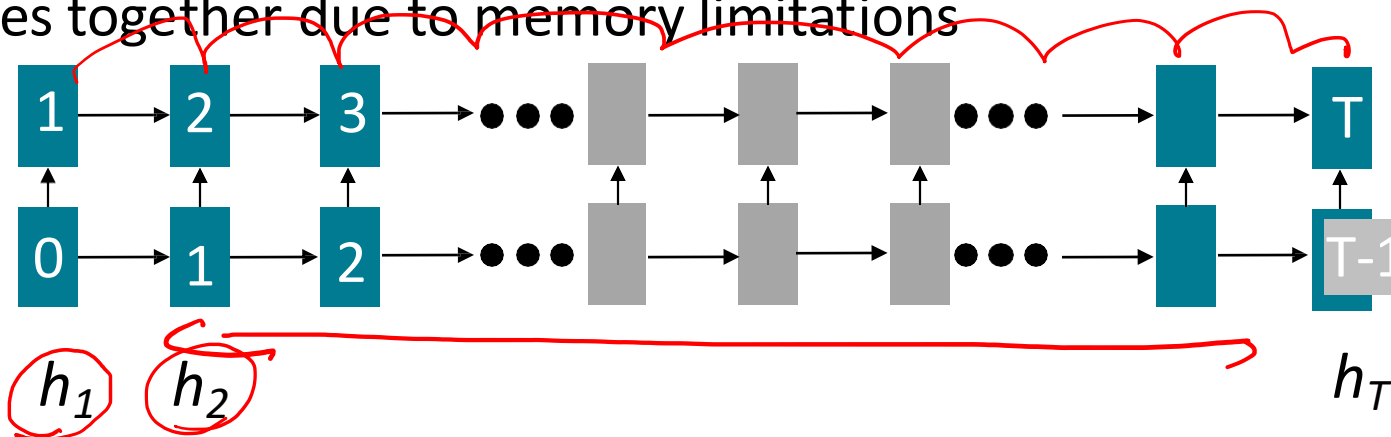**O(sequence length)** steps for distant word pairs to interact means:

- Hard to learn long-distance dependencies (because gradient problems!)
- Linear order of words is "baked in"; we already know sequential structure doesn't tell the whole story...



The **chef** who ...

*ate*

Info of **chef** has gone through O(sequence length) many layers!
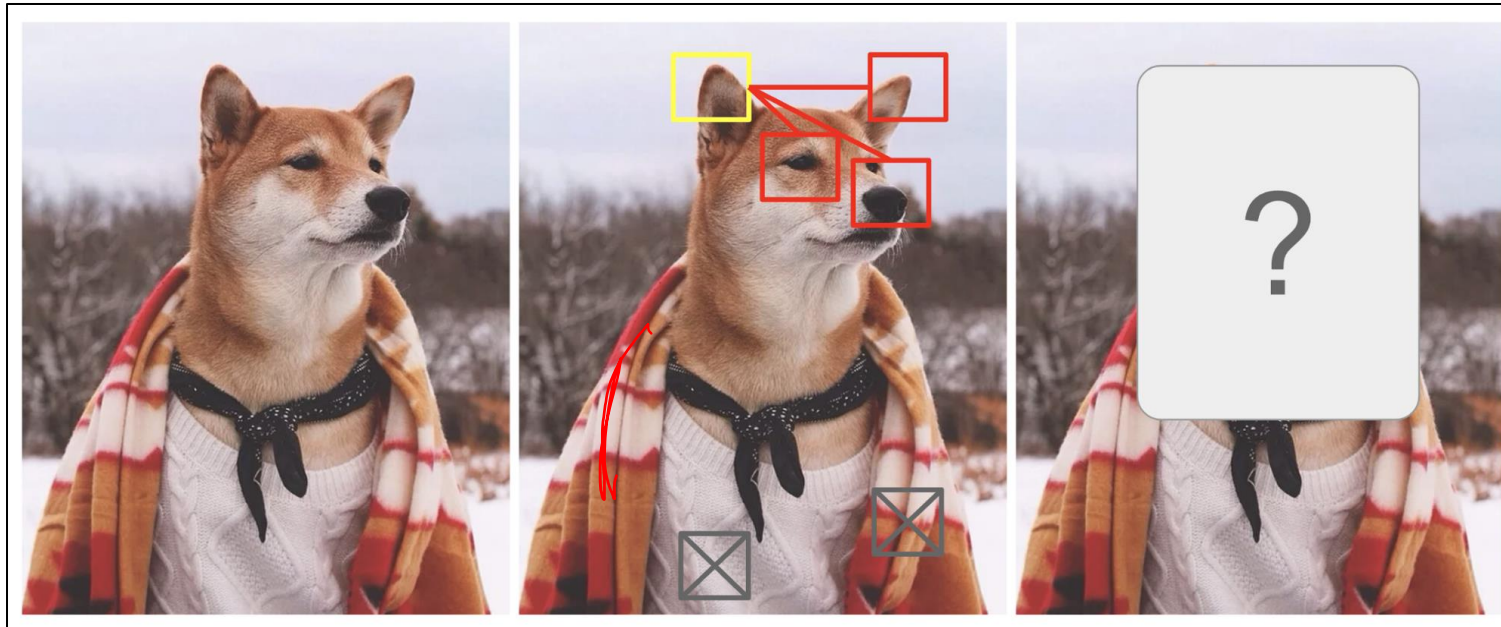
# Lack of  Parallelizability in RNN

- Forward and backward passes have **O(seq length)** unparallelizable operations
  - GPUs (and TPUs) can perform many independent computations at once!
  - But future RNN hidden states can't be computed in full before past RNN hidden states have been computed
  - Inhibits training on very large datasets!
  - Particularly problematic as sequence length increases, as we can no longer batch many examples together due to memory limitations



Numbers indicate min # of steps before a state can be computed

LET'S HAVE
A MOMENT
OF SILENCE
FOR ALL
THOSE
PEOPLE
DYING FOR
ATTENTION

# Attention



*A Shiba Inu in a men's outfit. The credit of the original photo goes to Instagram @mensweardog.*
*Source:* https://lilianweng.github.io/posts/2018-06-24-attention/

# Attention

# Attention Visuals/Animated

- Refer to separately uploaded Attention slides

"Noa can be annoying but she is a great cat.

cannot blindly use proximity

think of a way
to do this,
maybe automatically

"Noa can be annoying but she is a great cat.

cannot blindly use
proximity

"Noa can be annoying but she is a great cat.

think of a way to do this, maybe automatically

cannot blindly use proximity

$t_1$  $t_2$  $t_3$  $t_4$  $t_5$  $t_6$  $t_7$  $t_8$  $t_9$  $t_{10}$

$V_1$  $V_2$  $V_3$  $V_4$  $V_5$  $V_6$  $V_7$  $V_8$  $V_9$  $V_{10}$

$W_1$  $W_2$  $W_3$  $W_4$  $W_5$  $W_6$  $W_7$  $W_8$  $W_9$  $W_{10}$

looking for a reweighing method

$Y_1$  $Y_2$  $Y_3$  $Y_4$  $Y_5$  $Y_6$  $Y_7$  $Y_8$  $Y_9$  $Y_{10}$

more context

$V_k =$ [ ] k

$V_Q =$ [ ] k

family
royalty
history
power
gender
sentiment

son
daughter
cousin
land
country
army
own

dog
cat
running
swimming
help
nurse
computer

These words share meaning even if they're not in proximity!

Reweighing Plan

$$W_{kQ} = V_k \cdot V_Q \ ?$$

Let's explore this

Bank of the river.

↓ ↓ ↓ ↓

$V_1$   $V_2$   $V_3$   $V_4$



REWEIGH!        PROCESS

$y_1$   $y_2$   $y_3$   $y_4$   ← better, more context

$W_{11} V_1 + W_{12} V_2 + W_{13} V_3 + W_{14} V_4 = y_1$

$V_1 V_1 = W_{11}$

$V_1 V_2 = W_{12}$     $W_{11}$

$V_1 V_3 = W_{13}$     norm     $W_{12}$

$V_1 V_4 = W_{14}$              $W_{13}$

                               $W_{14}$

↳ sum to one

https://learning.rasa.com/transformers/self-attention/

17

Bank of the river.

↓ ↓ ↓ ↓

$V_1$    $V_2$    $V_3$    $V_4$



REWEIGH!     PROCESS

$y_1$    $y_2$    $y_3$    $y_4$    ← better, more context

$W_{11}V_1 + W_{12}V_2 + W_{13}V_3 + W_{14}V_4 = y_1$

↑      ↑      ↑      ↑

$V_1 \, V_1 = W_{11}$        $W_{11}$

$V_1 \, V_2 = W_{12}$    ⟶    $W_{12}$

$V_1 \, V_3 = W_{13}$    norm    $W_{13}$

$V_1 \, V_4 = W_{14}$        $W_{14}$

↑

⤷ sum to one

https://learning.rasa.com/transformers/self-attention/

18

Banh of the river.

$V_1 \quad V_2 \quad V_3 \quad V_4$

REWEIGH!    PROCESS

$Y_1 \quad Y_2 \quad Y_3 \quad Y_4$    ← better, more context

$V_1 V_1 = W_{11}$
$V_1 V_2 = W_{12}$     →    $W_{11}$
$V_1 V_3 = W_{13}$     norm    $W_{12}$
$V_1 V_4 = W_{14}$            $W_{13}$
                              $W_{14}$
            ↳ sum to one        ↑

$W_{11}V_1 + W_{12}V_2 + W_{13}V_3 + W_{14}V_4 = Y_1$

reweigh all vectors towards $V_1$

Bank of the river.

$$W_{11}V_1 + W_{12}V_2 + W_{13}V_3 + W_{14}V_4 = Y_1$$

$$W_{21}V_1 + W_{22}V + W_{23}V_3 + W_{24}V_4 = Y_2$$

$$W_{31}V_1 + W_{32}V_2 + W_{33}V_3 + W_{34}V_4 = Y_3$$

$$W_{41}V_1 + W_{42}V_2 \rightarrow W_{43}V_3 + W_{44}V_4 = Y_4$$

$V_1 \quad V_2 \quad V_3 \quad V_4$

REWEIGH!    PROCESS

$Y_1 \quad Y_2 \quad Y_3 \quad Y_4$    ← better, more context

$$V_1 V_1 = W_{11}$$
$$V_1 V_2 = W_{12}$$
$$V_1 V_3 = W_{13}$$
$$V_1 V_4 = W_{14}$$

norm →

$W_{11}$
$W_{12}$
$W_{13}$
$W_{14}$

— sum to one

Banh of the river.

$V_1 \quad V_2 \quad V_3 \quad V_4$

REWEIGH!

PROCESS

$y_1 \quad y_2 \quad y_3 \quad y_4$

← better, more context

$$W_{11}V_1 + W_{12}V_2 + W_{13}V_3 + W_{14}V_4 = y_1$$

$$W_{21}V_1 + W_{22}V + W_{23}V_3 + W_{24}V_4 = y_2$$

$$W_{31}V_1 + W_{32}V_2 + W_{33}V_3 + W_{34}V_4 = y_3$$

$$W_{41}V_1 + W_{42}V_2 + W_{43}V_3 + W_{44}V_4 = y_4$$

- I've not trained any weights
- Order has no influence
- Proximity has no influence
- Shape independant

$$V_1 V_1 = W_{11}$$
$$V_1 V_2 = W_{12}$$
$$V_1 V_3 = W_{13}$$
$$V_1 V_4 = W_{14}$$

norm →

$W_{11}$
$W_{12}$
$W_{13}$
$W_{14}$

sum to one

SELF ATTENTION ←

https://learning.rasa.com/transformers/self-attention/

$V_1 \quad V_2 \quad V_3 \quad V_4 \quad\quad\quad\quad \longrightarrow y_4$

$\longrightarrow y_3$

$W_{31} \quad W_{32} \quad W_{33} \quad W_{34} \quad\quad\quad\quad \longrightarrow y_2$

$\longrightarrow y_1$

NORMALISE

$$\sum_j W_{3j} = 1$$

$S_{31} \quad S_{32} \quad S_{33} \quad S_{34}$

$y_1 \quad y_2 \quad y_3 \quad y_4$

no weights $\longrightarrow$ ATTENTION

DOT PRODUCT $\leftarrow V_3$

$V_1 \quad V_2 \quad V_3 \quad V_4$

$V_1 \quad V_2 \quad V_3 \quad V_4$

23

$V_1 \quad V_2 \quad V_3 \quad V_4$

$\times$

$Y_4$

$Y_3$

$+$

$Y_2$

$Y_1$

$W_{31} \quad W_{32} \quad W_{33} \quad W_{34}$

NORMALISE

$\sum_j W_{3j} = 1$

$S_{31} \quad S_{32} \quad S_{33} \quad S_{34}$

DOT PRODUCT

$\leftarrow V_3$

~~no~~ weights $\longrightarrow$

$V_1 \quad V_2 \quad V_3 \quad V_4$

$Y_1 \quad Y_2 \quad Y_3 \quad Y_4$

ATTENTION

$V_1 \quad V_2 \quad V_3 \quad V_4$

24

$V_1 \quad V_2 \quad V_3 \quad V_4$

$\times$

$\times$

$\times$

$\times$

$+ \quad \longrightarrow Y_3$

$\longrightarrow Y_4$

$\longrightarrow Y_2$

$\longrightarrow Y_1$

$W_{31} \quad W_{32} \quad W_{33} \quad W_{34}$

NORMALISE

$\sum_j W_{3j} = 1$

$S_{31} \quad S_{32} \quad S_{33} \quad S_{34}$

DOT PRODUCT

$\leftarrow V_3$

$V_1 \quad V_2 \quad V_3 \quad V_4$

26

$$\sum_j W_{3j} = 1$$

NORMALISE

$S_{31} \quad S_{32} \quad S_{33} \quad S_{34}$

DOT PRODUCT $\leftarrow V_3$

$V_1 \quad V_2 \quad V_3 \quad V_4$

$V_1 \quad V_2 \quad V_3 \quad V_4$

$W_{31} \quad W_{32} \quad W_{33} \quad W_{34}$

$Y_4 \quad Y_3 \quad Y_2 \quad Y_1$

$$V_1 \qquad V_2 \qquad V_3 \qquad V_4 \qquad\qquad\qquad \to Y_4$$

$$\to Y_3$$

$$\times$$

$$+ \to Y_3$$

$$\to Y_2$$

$$\to Y_1$$

$$W_{31} \qquad W_{32} \qquad W_{33} \qquad W_{34}$$

NORMALISE

$$\sum_j W_{3j} = 1$$

$$S_{31} \qquad S_{32} \qquad S_{33} \qquad S_{34}$$

QUERY

DOT PRODUCT

$$\leftarrow V_3$$

$$V_1 \qquad V_2 \qquad V_3 \qquad V_4$$

28

$$\sum_j w_{3j} = 1$$

NORMALISE

$S_{31} \quad S_{32} \quad S_{33} \quad S_{34}$

DOT PRODUCT $\leftarrow V_{\underline{3}}$ — QUERY

$V_1 \quad V_2 \quad V_3 \quad V_4$ $\longleftarrow$ KEYS

VALUES

$V_1$   $V_2$   $V_3$   $V_4$

$\times$

$\times$

$\times$

$\times$

$+$ → $Y_4$

→ $Y_3$

→ $Y_2$

→ $Y_1$

$W_{31}$   $W_{32}$   $W_{33}$   $W_{34}$

$\sum_j W_{3j} = 1$

NORMALISE

$S_{31}$   $S_{32}$   $S_{33}$   $S_{34}$

QUERY

DOT PRODUCT

← $V_3$

KEYS

$V_1$   $V_2$   $V_3$   $V_4$

# Key, Value and Query
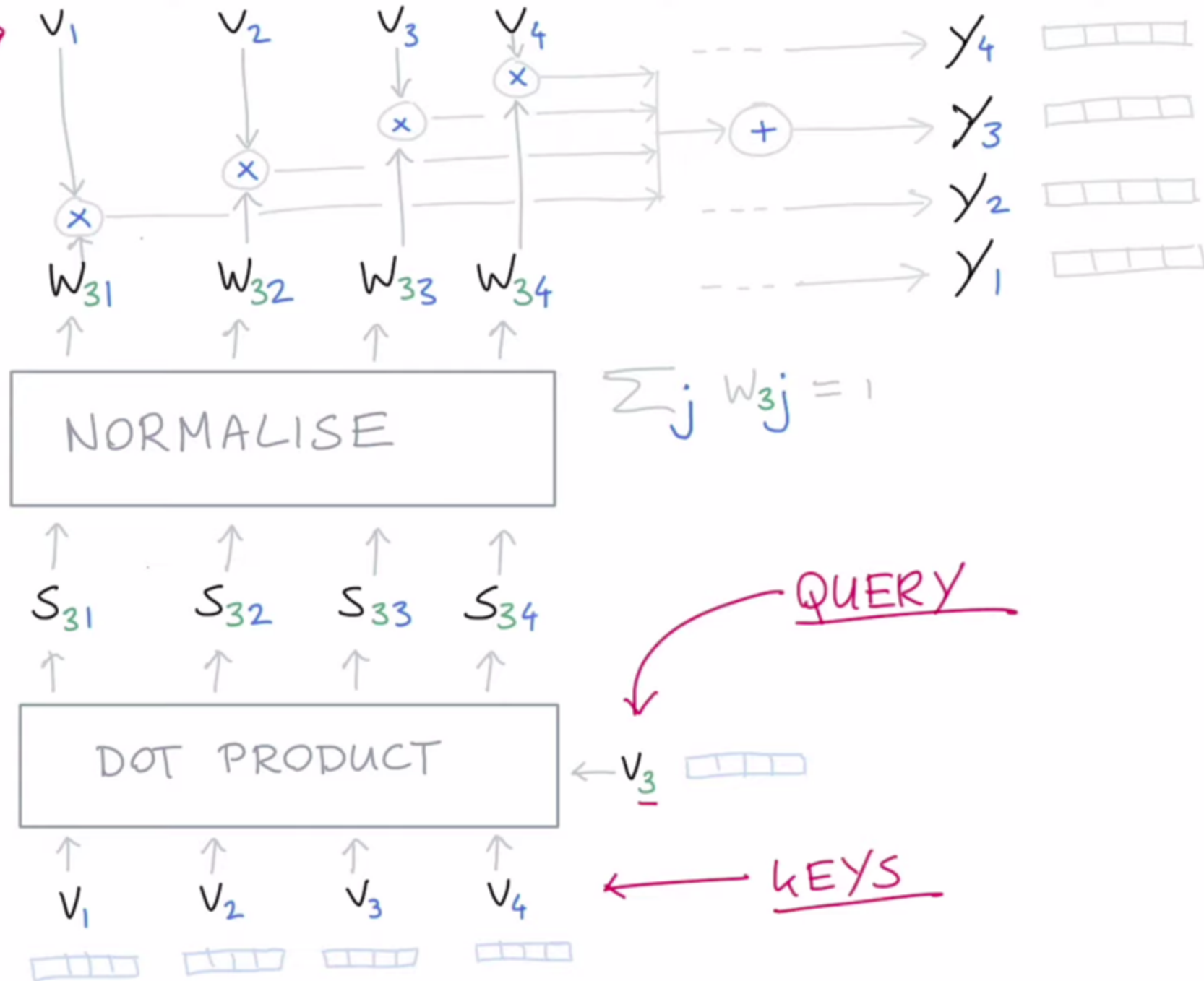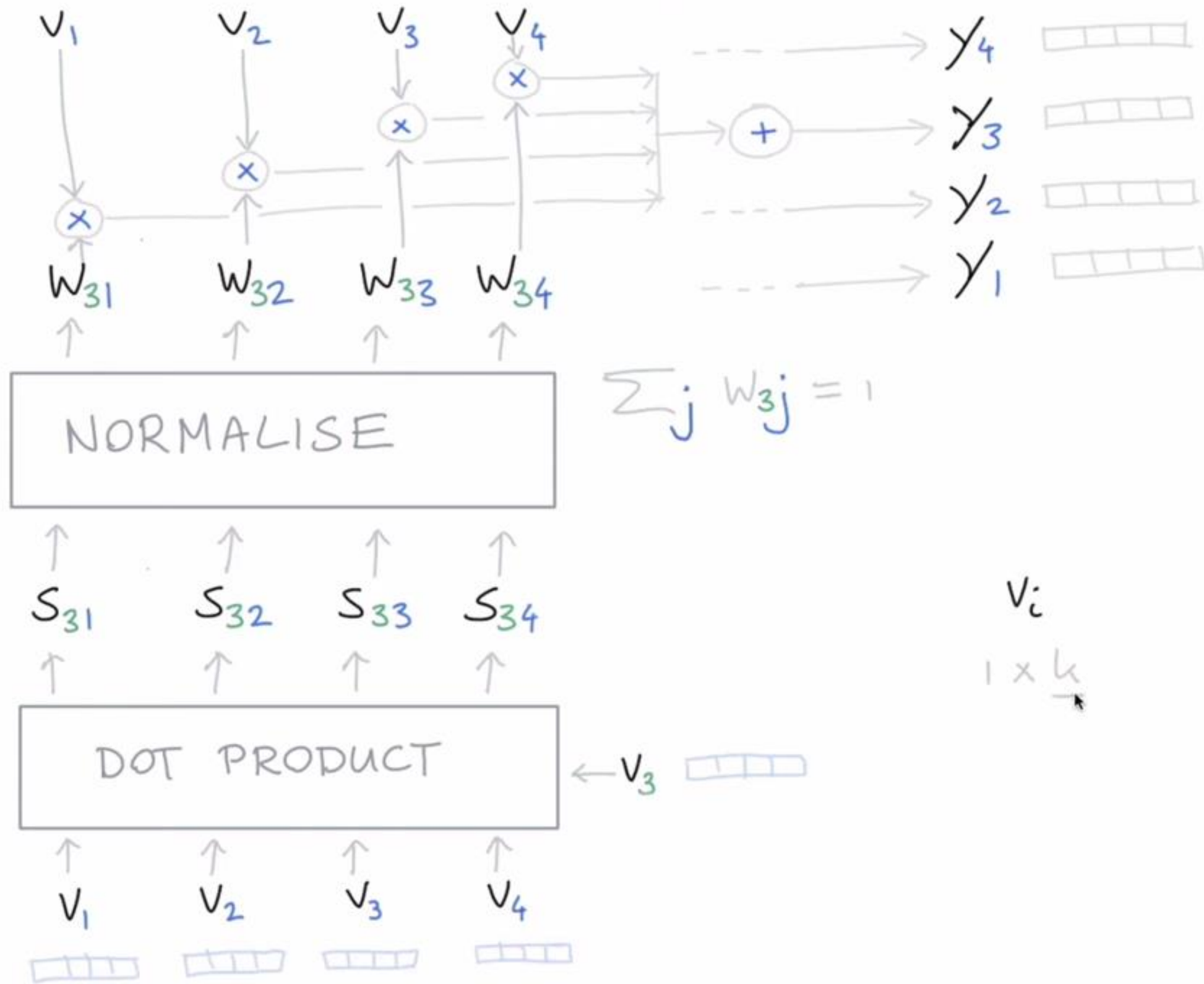
- The major component in the transformer is the unit of <span style="color:red">multi-head self-attention mechanism.</span>

$V_1 \quad V_2 \quad V_3 \quad V_4 \longrightarrow Y_4$

$\longrightarrow Y_3$

$W_{31} \quad W_{32} \quad W_{33} \quad W_{34} \longrightarrow Y_2$

$\longrightarrow Y_1$

NORMALISE

$\sum_j W_{3j} = 1$

$S_{31} \quad S_{32} \quad S_{33} \quad S_{34}$

$V_i$

$1 \times k$

DOT PRODUCT $\leftarrow V_3$

$V_1 \quad V_2 \quad V_3 \quad V_4$

32

$V_1$  $V_2$  $V_3$  $V_4$  → $Y_4$

$\times$  → $Y_3$

$\times$  $\oplus$

$\times$

$\times$  → $Y_2$

$W_{31}$  $W_{32}$  $W_{33}$  $W_{34}$  → $Y_1$

$\sum_j W_{3j} = 1$

NORMALISE

$S_{31}$  $S_{32}$  $S_{33}$  $S_{34}$

$V_i$  $M$

$1 \times \underline{k}$

DOT PRODUCT  ← $V_3$

$V_1$  $V_2$  $V_3$  $V_4$

33

$V_1 \quad V_2 \quad V_3 \quad V_4 \qquad\qquad\qquad \longrightarrow Y_4$

$\qquad\qquad\qquad\qquad\qquad\qquad \longrightarrow Y_3$

$\qquad\qquad\qquad\qquad\qquad\qquad +\longrightarrow Y_2$

$W_{31} \quad W_{32} \quad W_{33} \quad W_{34} \qquad\qquad \longrightarrow Y_1$

NORMALISE

$\sum_j W_{3j} = 1$

$S_{31} \quad S_{32} \quad S_{33} \quad S_{34}$
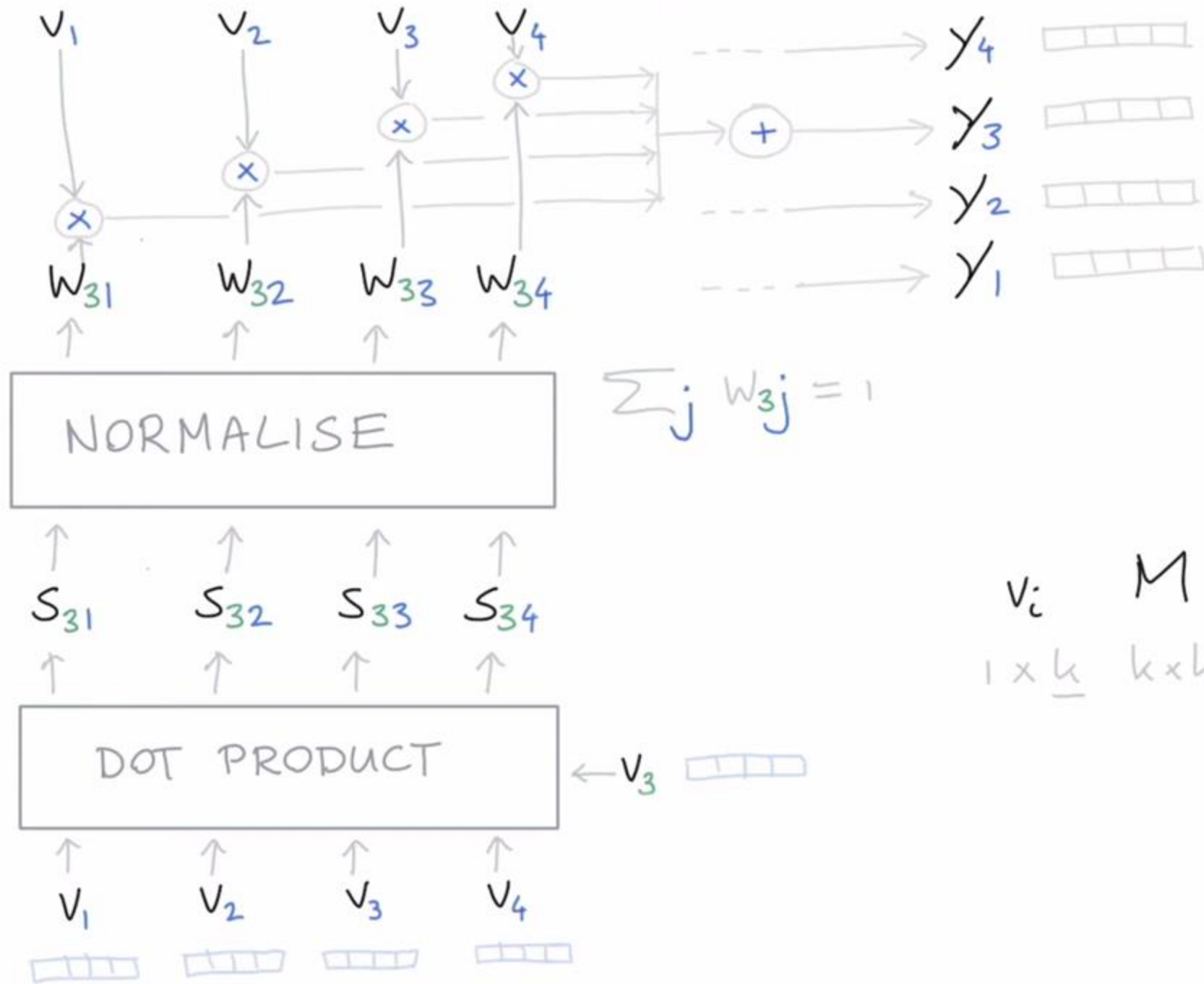
$V_i \qquad M$

$1 \times \underline{k} \quad k \times k$

DOT PRODUCT $\quad \leftarrow V_3$

$V_1 \quad V_2 \quad V_3 \quad V_4$

$V_1$  $V_2$  $V_3$  $V_4$

$Y_4$

$Y_3$

$Y_2$

$Y_1$

$W_{31}$  $W_{32}$  $W_{33}$  $W_{34}$

$$\sum_j W_{3j} = 1$$

NORMALISE

$S_{31}$  $S_{32}$  $S_{33}$  $S_{34}$

$V_i$  $M = [\ ]$

$1 \times \underline{k}$  $k \times k$

DOT PRODUCT

$\leftarrow V_3$

$V_1$  $V_2$  $V_3$  $V_4$

35

$v_1 \quad v_2 \quad v_3 \quad v_4 \qquad\qquad y_4$

$y_3$

$y_2$

$y_1$

$W_{31} \quad W_{32} \quad W_{33} \quad W_{34}$

$\sum_j W_{3j} = 1$

NORMALISE

$S_{31} \quad S_{32} \quad S_{33} \quad S_{34}$

$v_i \quad M = [\ ]$

$1 \times \underline{k} \quad k \times k \quad 1 \times k$

DOT PRODUCT

$\leftarrow v_3$

$v_1 \quad v_2 \quad v_3 \quad v_4$

$V_1$    $V_2$    $V_3$    $V_4$

$Y_4$

$Y_3$

$Y_2$

$Y_1$

$W_{31}$   $W_{32}$   $W_{33}$   $W_{34}$

$\sum_j W_{3j} = 1$

NORMALISE

$S_{31}$   $S_{32}$   $S_{33}$   $S_{34}$

$V_i$   $M = [\ ]$

$1 \times \underline{k}$   $k \times k$   $1 \times k$

DOT PRODUCT   $\leftarrow V_3$

$V_1 M_k$   $V_2 M_k$   $V_3 M_k$   $V_4 M_k$   KEY

37

$V_1$   $V_2$   $V_3$   $V_4$ → $y_4$

→ $y_3$

→ $y_2$

$W_{31}$   $W_{32}$   $W_{33}$   $W_{34}$ → $y_1$

$\sum_j W_{3j} = 1$

NORMALISE

$S_{31}$   $S_{32}$   $S_{33}$   $S_{34}$

$V_i$   $M = []$

$1 \times k$   $k \times k$   $1 \times k$

DOT PRODUCT   ← $V_3$

$V_1 M_k$   $V_2 M_k$   $V_3 M_k$   $V_4 M_k$   KEY

38

VALUE

$v_1 M_v$    $v_2 M_v$    $v_3 M_v$    $v_4 M_v$    --- → $y_4$   ☐☐☐☐

× → $y_3$   ☐☐☐☐

× → + → $y_3$

× → $y_2$   ☐☐☐☐☐

×

$W_{31}$    $W_{32}$    $W_{33}$    $W_{34}$    --- → $y_1$   ☐☐☐☐

↑        ↑         ↑         ↑

$$\boxed{\text{NORMALISE}}$$    $\sum_j W_{3j} = 1$

↑        ↑         ↑         ↑

$S_{31}$    $S_{32}$    $S_{33}$    $S_{34}$        QUERY        $v_i$   $M = []$

↑        ↑         ↑         ↑                              $1 \times k$   $k \times k$   $1 \times k$

$$\boxed{\text{DOT PRODUCT}}$$    ← $v_3 M_Q$ ☐☐☐

↑        ↑         ↑         ↑

$v_1 M_k$    $v_2 M_k$    $v_3 M_k$    $v_4 M_k$    KEY

☐☐☐    ☐☐☐☐    ☐☐☐    ☐☐☐

SELF ATTENTION

KEYS — LINEAR

QUERIES — LINEAR

VALUES — LINEAR

$\begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}$

$v_i$

MATMUL → $[S_{ij}]$

NORMALISE → $[W_{ij}]$

MATMUL

$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$

SELF ATTENTION

KEYS
LINEAR

QUERIES
LINEAR

VALUES
LINEAR

$\begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}$

$v_i$ ▭▭▭

MATMUL

$[S_{ij}]$

NORMALISE

$\rightarrow [W_{ij}]$

MATMUL

$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$

SELF ATTENTION

SELF ATTENTION

$V_1 \qquad V_2 \qquad V_3$

SELF ATTENTION

KEYS
LINEAR

QUERIES
LINEAR

VALUES
LINEAR

MATMUL

$[S_{ij}]$

NORMALISE

$[W_{ij}]$

MATMUL

$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$

$\begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}$

$V_i$ ▭▭▭

NER

SELF ATTENTION

SELF ATTENTION

$V_1$    $V_2$    $V_3$

43

SELF ATTENTION

KEYS
LINEAR

QUERIES
LINEAR

VALUES
LINEAR

MATMUL

$[S_{ij}]$

NORMALISE

$\rightarrow [W_{ij}]$

MATMUL

$\begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}$

$V_i$ ▭▭▭

$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$

NER

SELF ATTENTION

SELF ATTENTION

$V_1$   $V_2$   $V_3$

44

SELF ATTENTION

KEYS
LINEAR

QUERIES
LINEAR

VALUES
LINEAR

$\begin{bmatrix} V_1 \\ \vdots \\ V_n \end{bmatrix}$

$V_i$ ▯▯▯▯

MATMUL

$[S_{ij}]$

NORMALISE → $[W_{ij}]$

MATMUL

$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$

NER

SELF ATTENTION

SELF ATTENTION

$V_1 \quad V_2 \quad V_3$

SELF ATTENTION

KEYS

LINEAR

QUERIES

LINEAR

VALUES

LINEAR

$\begin{bmatrix} V_1 \\ \vdots \\ V_n \end{bmatrix}$

$V_i$

MATMUL

$[S_{ij}]$

NORMALISE

$[W_{ij}]$

MATMUL

$\begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}$

NER

SELF ATTENTION

SELF ATTENTION

$V_1$          $V_2$          $V_3$

46

SELF ATTENTION

KEYS

LINEAR

QUERIES

LINEAR

VALUES

LINEAR

MATMUL

$[S_{ij}]$

NORMALISE

$[W_{ij}]$

MATMUL

$\begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}$

$v_i$ ▢▢▢

$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$

NER

SELF ATTENTION

SELF ATTENTION

$v_1$   $v_2$   $v_3$

47

SELF ATTENTION

kEYS
LINEAR

QUERIES
LINEAR

VALUES
LINEAR

MATMUL

$[S_{ij}]$

NORMALISE

$[W_{ij}]$

MATMUL

$\begin{bmatrix} V_1 \\ \vdots \\ V_n \end{bmatrix}$

$V_i$

$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$

NER

SELF ATTENTION

SELF ATTENTION

$V_1$    $V_2$    $V_3$

https://learning.rasa.com/transformers/self-attention/

SELF ATTENTION

KEYS
LINEAR

QUERIES
LINEAR

VALUES
LINEAR

$\begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}$

$V_i$

MATMUL

$[S_{ij}]$

NORMALISE

$[W_{ij}]$

MATMUL

$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$

NER

SELF ATTENTION

SELF ATTENTION

$V_1$   $V_2$   $V_3$

49

SELF ATTENTION

kEYS
LINEAR

QUERIES
LINEAR

VALUES
LINEAR

$V_i$

MATMUL

$[S_{ij}]$

NORMALISE

$[W_{ij}]$

MATMUL

$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$

$\begin{bmatrix} V_1 \\ \vdots \\ V_n \end{bmatrix}$
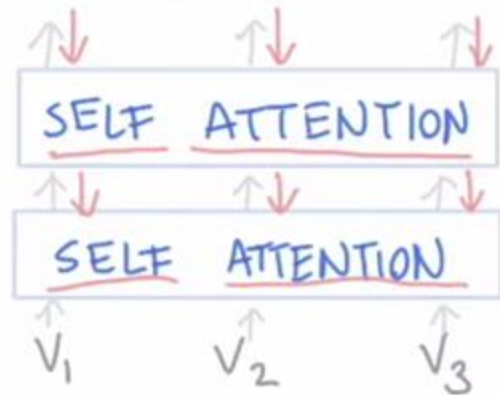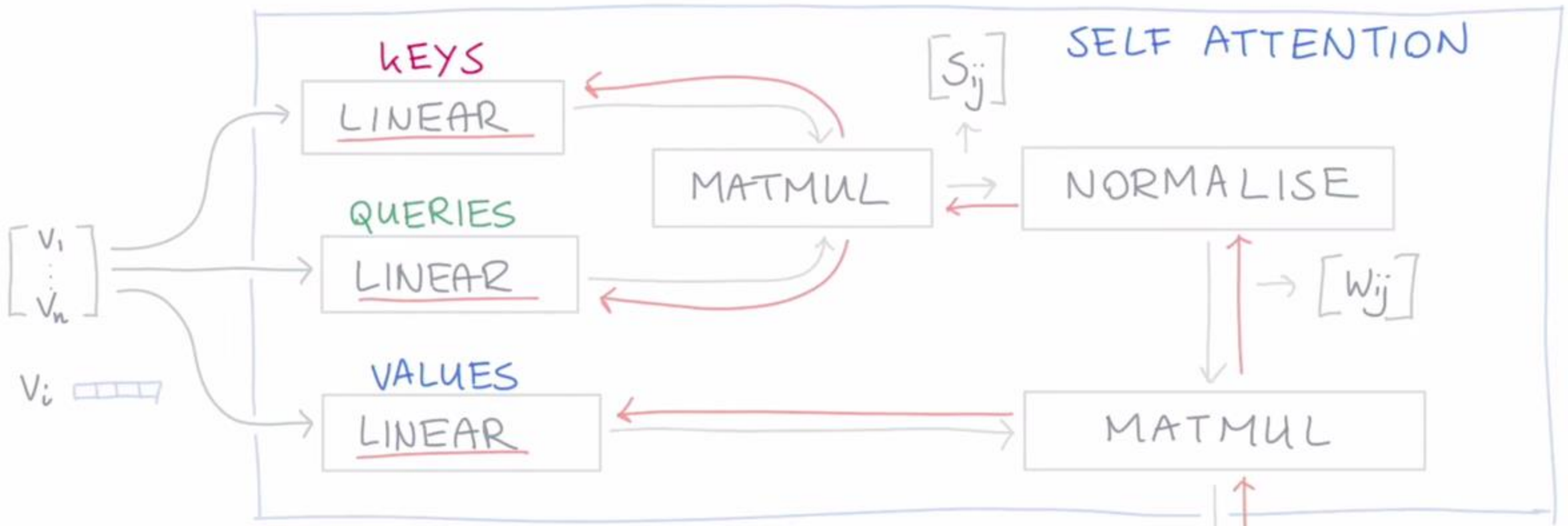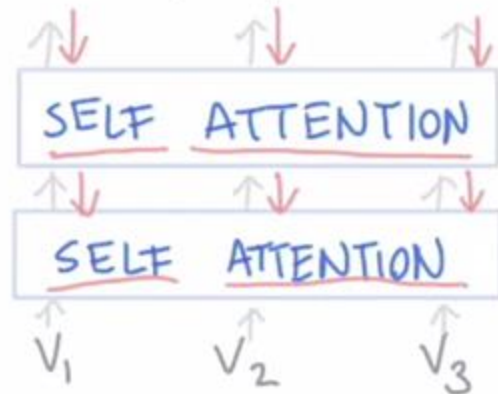
NER

SELF ATTENTION

SELF ATTENTION

$V_1$  $V_2$  $V_3$

← clich together ⸭

# Acknowledgments

These slides were adapted from the book

SPEECH and LANGUAGE PROCESSING: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition

and

some modifications from presentations and resources found in the WEB by several scholars mentioned in references.

# References

- https://slds-lmu.github.io/seminar_nlp_ss20/attention-and-self-attention-for-nlp.html
- Attention? Attention! | Lil'Log (lilianweng.github.io)

# Reference materials

- [https://vlanc-lab.github.io/mu-nlp-course/teachings/fall2024-AI-schedule.html](https://vlanc-lab.github.io/mu-nlp-course/teachings/fall2024-AI-schedule.html)

- Lecture notes
- (A) Speech and Language Processing by Daniel Jurafsky and James H. Martin
- (B) Natural Language Processing with Python. (updated edition based on Python 3 and NLTK 3) Steven Bird et al. O'Reilly Media