# Natural Language Processing

## CS 3216/UG, AI 5203/PG

### Week-8
### 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
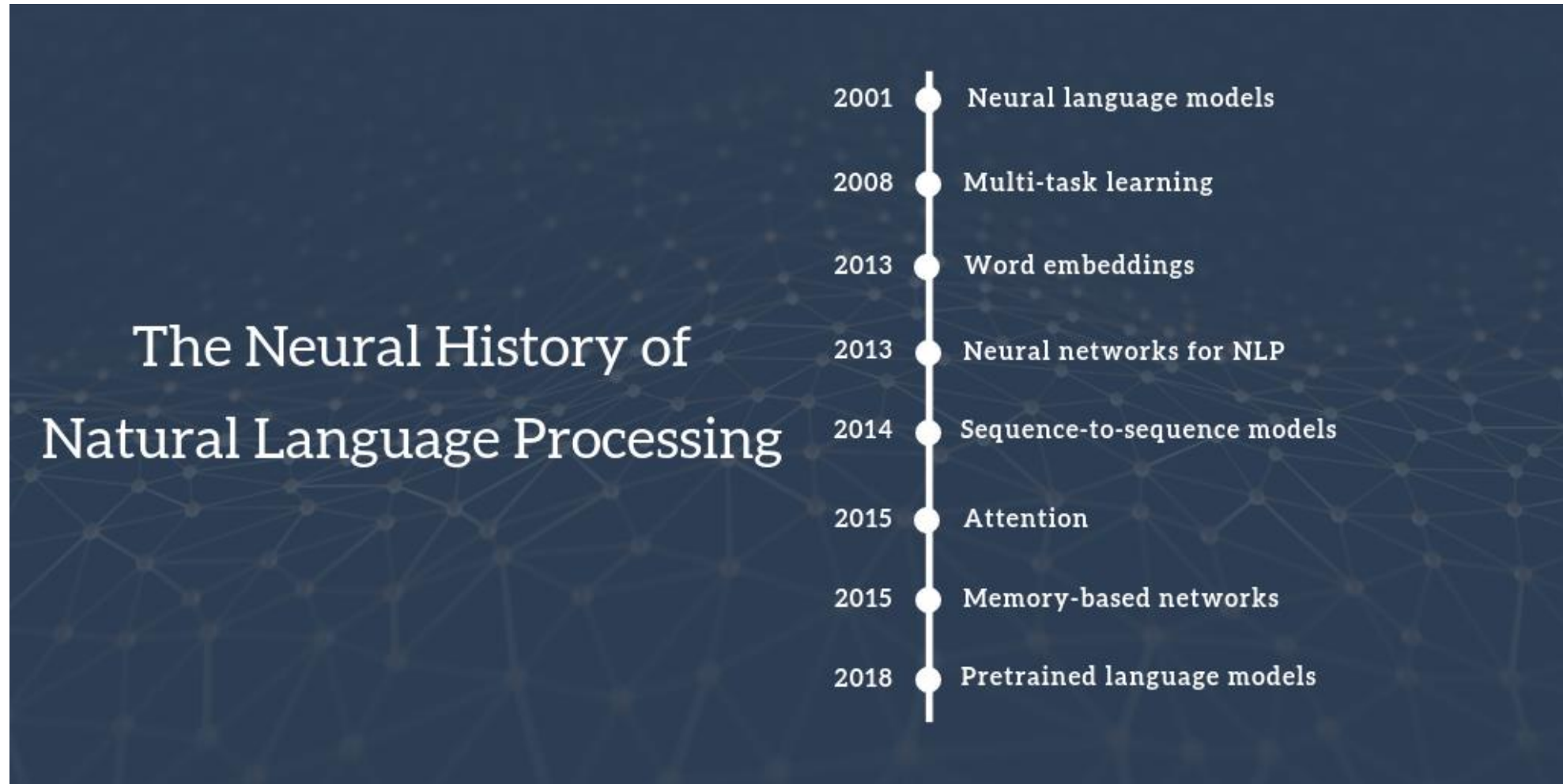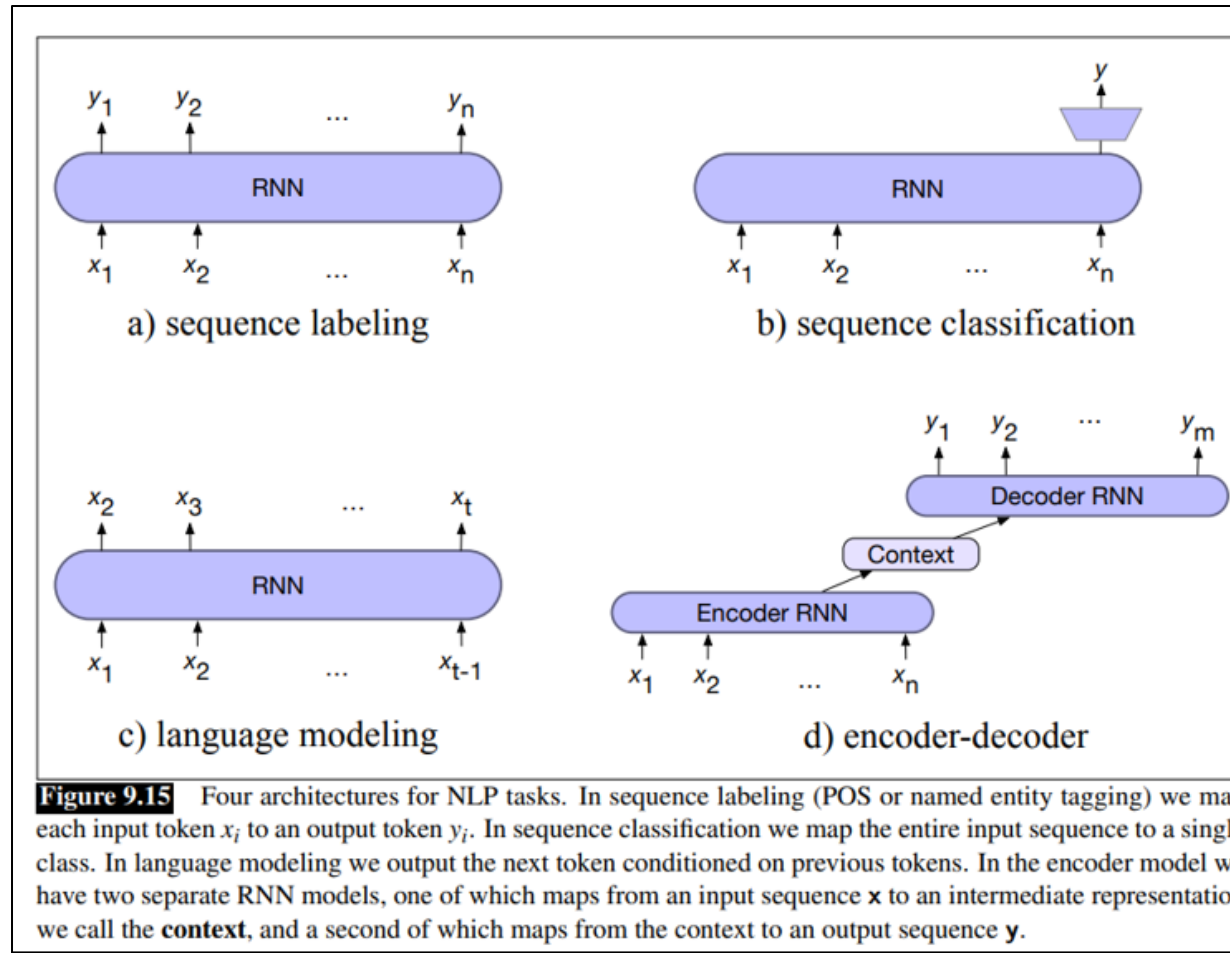
# History of Neural models in NLP



The Neural History of Natural Language Processing

| | |
|---|---|
| 2001 | Neural language models |
| 2008 | Multi-task learning |
| 2013 | Word embeddings |
| 2013 | Neural networks for NLP |
| 2014 | Sequence-to-sequence models |
| 2015 | Attention |
| 2015 | Memory-based networks |
| 2018 | Pretrained language models |

Image source : https://www.ruder.io/a-review-of-the-recent-history-of-nlp/

# Different variants of RNN

- Stacked RNN
- Bi-directional RNN
- Many more

# Sequence -to-Sequence learning



**Figure 9.15** Four architectures for NLP tasks. In sequence labeling (POS or named entity tagging) we map each input token $x_i$ to an output token $y_i$. In sequence classification we map the entire input sequence to a single class. In language modeling we output the next token conditioned on previous tokens. In the encoder model we have two separate RNN models, one of which maps from an input sequence **x** to an intermediate representation we call the **context**, and a second of which maps from the context to an output sequence **y**.

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

# Sequence to Sequence Learning with Neural Networks

**Ilya Sutskever**
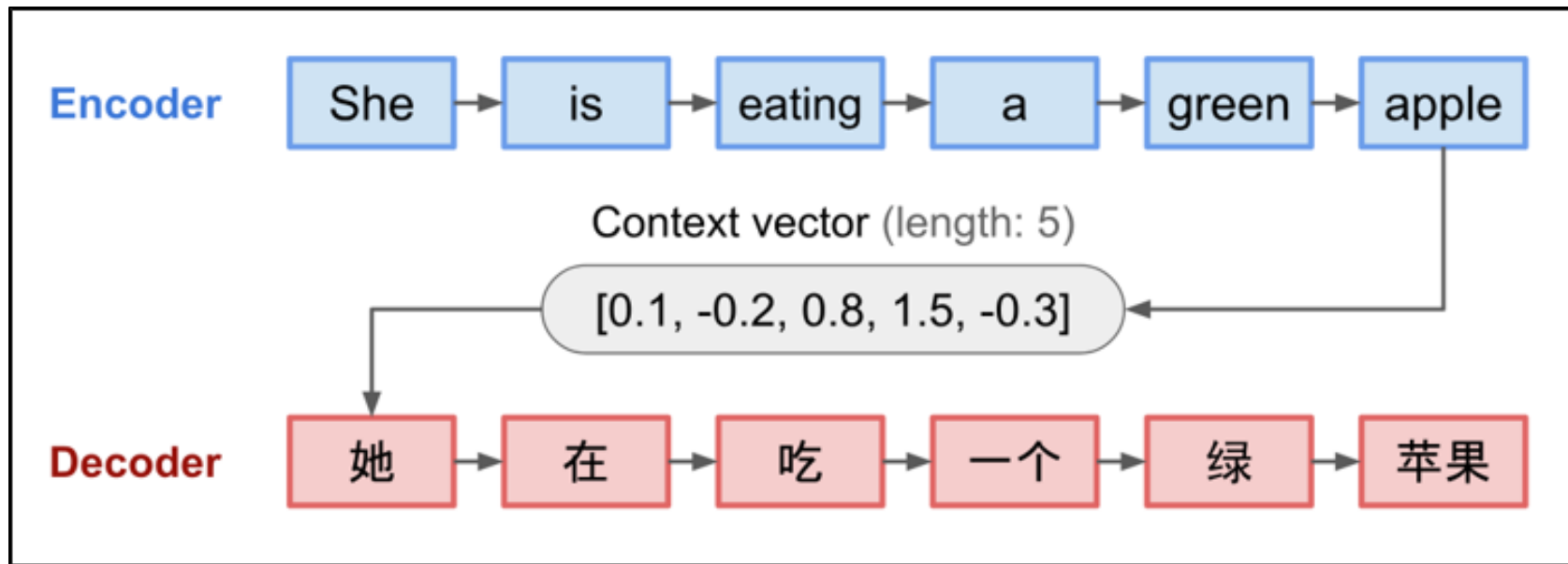Google
ilyasu@google.com

**Oriol Vinyals**
Google
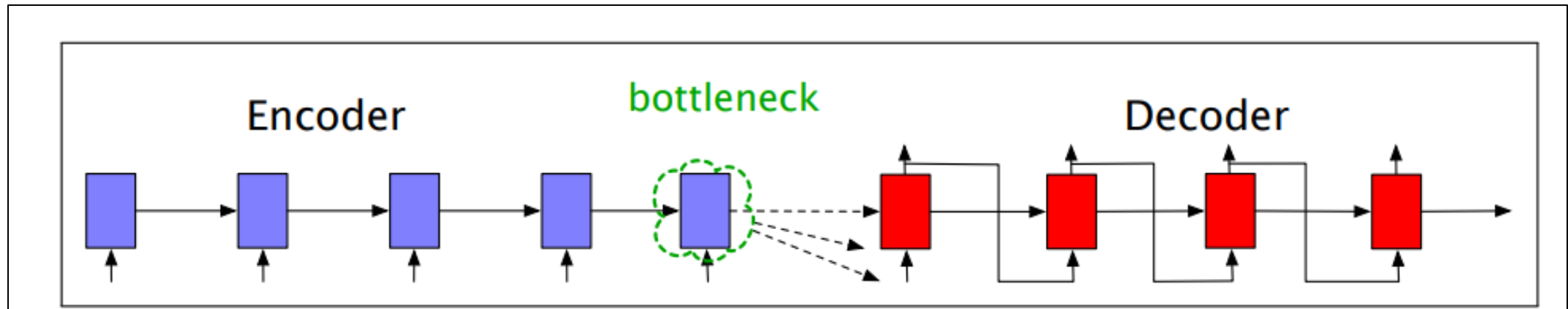vinyals@google.com

**Quoc V. Le**
Google
qvl@google.com

## Abstract

Deep Neural Networks (DNNs) are powerful models that have achieved excellent performance on difficult learning tasks. Although DNNs work well whenever large labeled training sets are available, they cannot be used to map sequences to sequences. In this paper, we present a general end-to-end approach to sequence learning that makes minimal assumptions on the sequence structure. Our method uses a multilayered Long Short-Term Memory (LSTM) to map the input sequence to a vector of a fixed dimensionality, and then another deep LSTM to decode the target sequence from the vector. Our main result is that on an English to French translation task from the WMT'14 dataset, the translations produced by the LSTM achieve a BLEU score of 34.8 on the entire test set, where the LSTM's BLEU score was penalized on out-of-vocabulary words. Additionally, the LSTM did not have difficulty on long sentences. For comparison, a phrase-based SMT system achieves a BLEU score of 33.3 on the same dataset. When we used the LSTM to rerank the 1000 hypotheses produced by the aforementioned SMT system, its BLEU score increases to 36.5, which is close to the previous best result on this task. The LSTM also learned sensible phrase and sentence representations that are sensitive to word order and are relatively invariant to the active and the pas-

https://arxiv.org/pdf/1409.3215

6

# 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
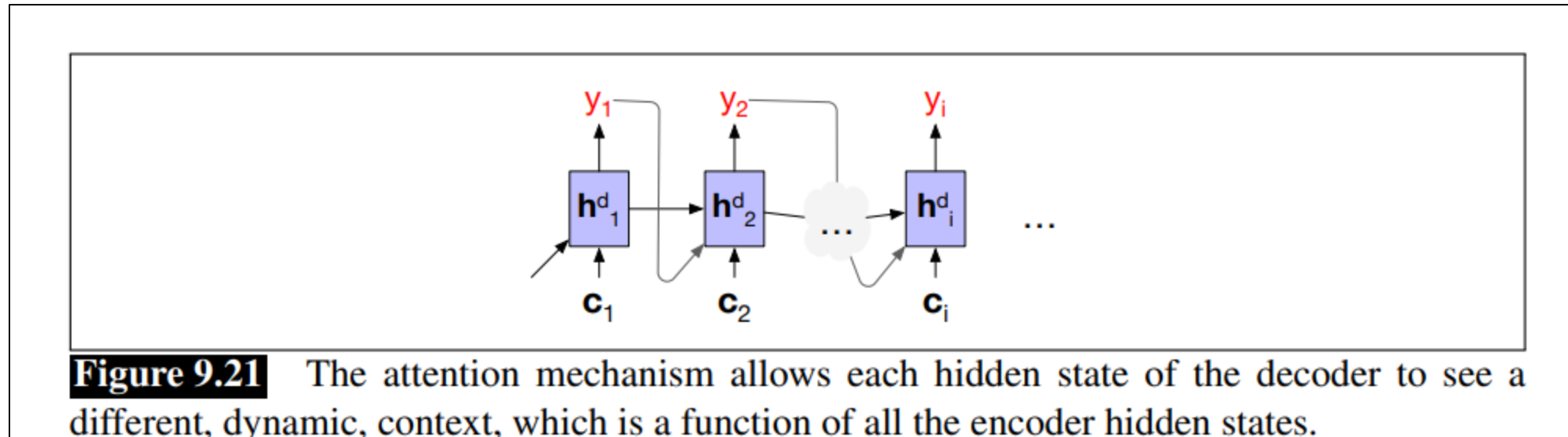  - incapability of remembering long sentences

**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.

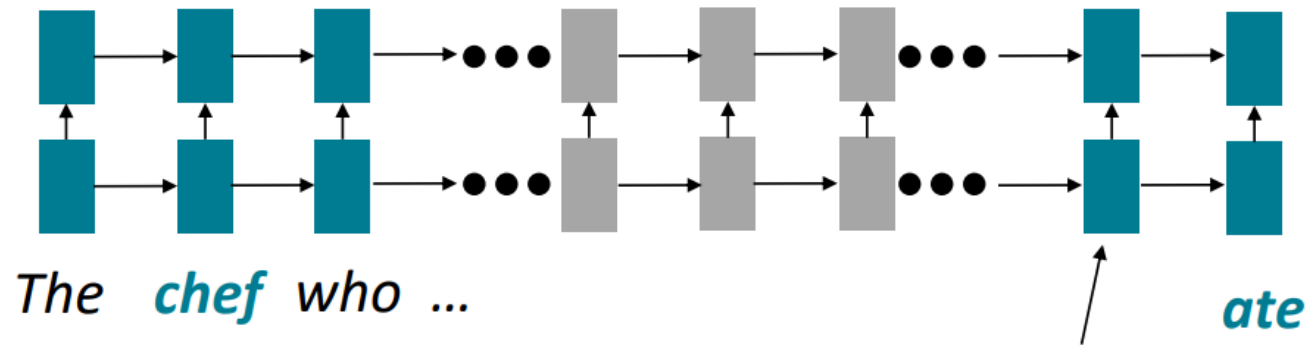https://lena-voita.github.io/nlp_course/seq2seq_and_attention.html

# 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
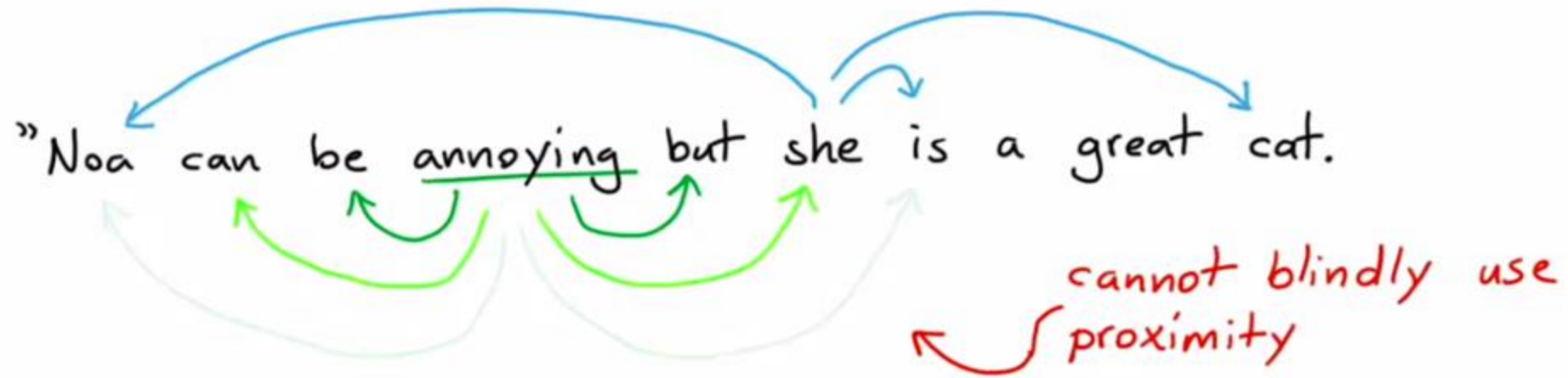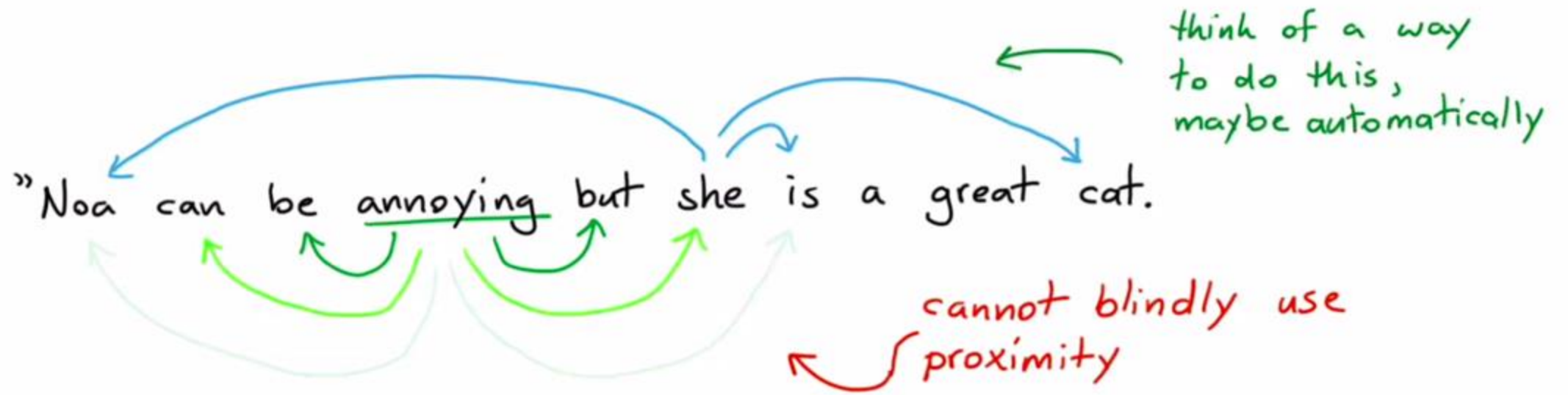PEOPLE
DYING FOR
ATTENTION

# Attention



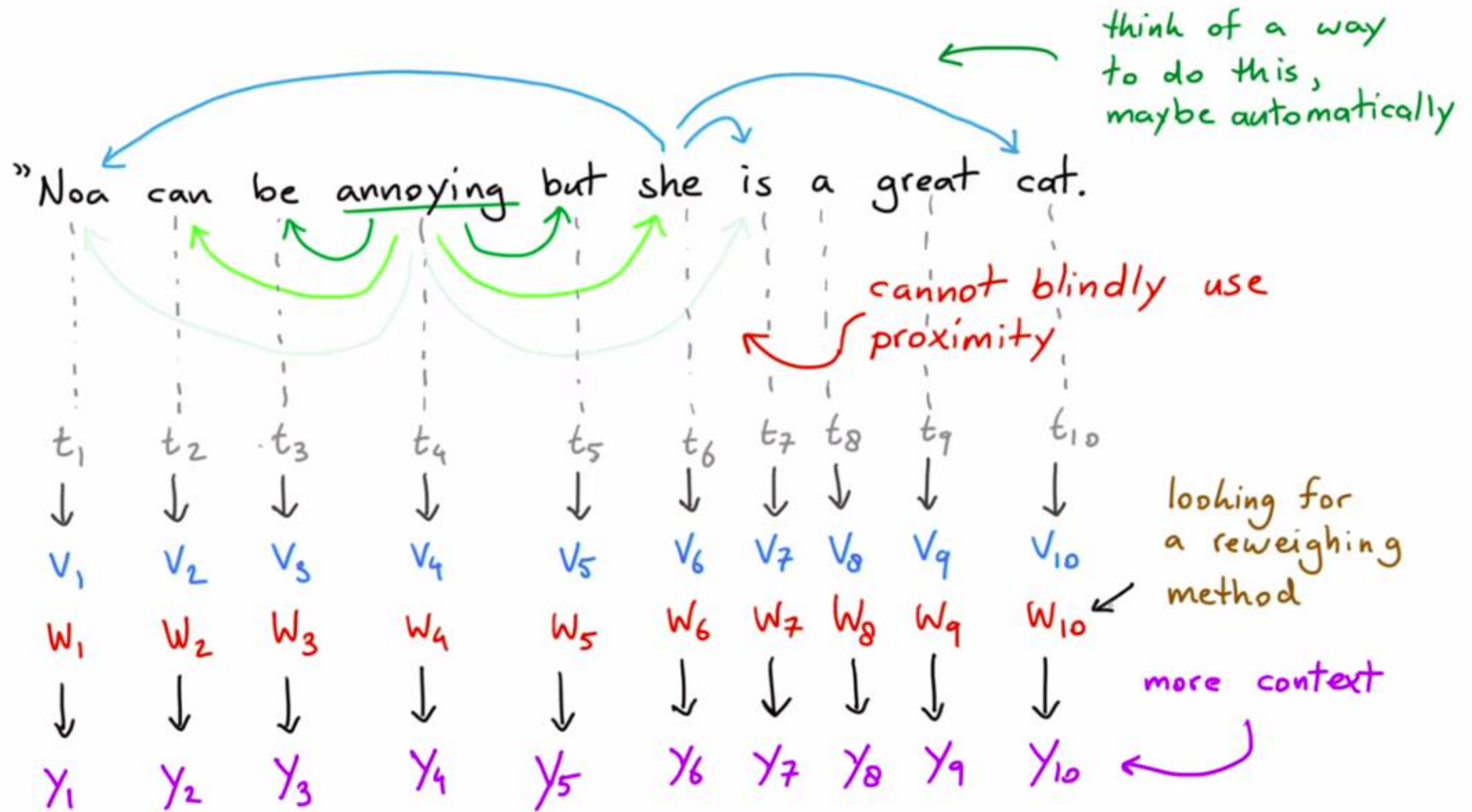*A Shiba Inu in a men's outfit. The credit of the original photo goes to Instagram @mensweardog.*
*Source: https://lilianweng.github.io/posts/2018-06-24-attention/*

# Attention

"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

think of a way to do this, maybe automatically

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

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

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

$V_k =$ [boxes with colored cells] $k$

$V_Q =$ [boxes with colored cells] $k$

family
royalty
history
power
gender
sentiment

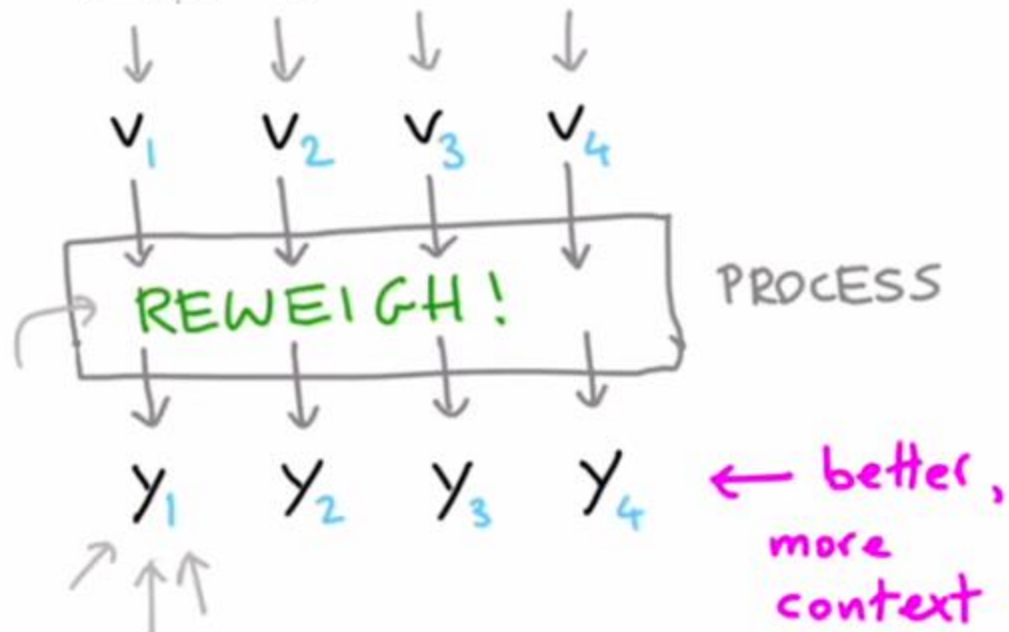**Reweighing Plan**

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

Let's explore this

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!

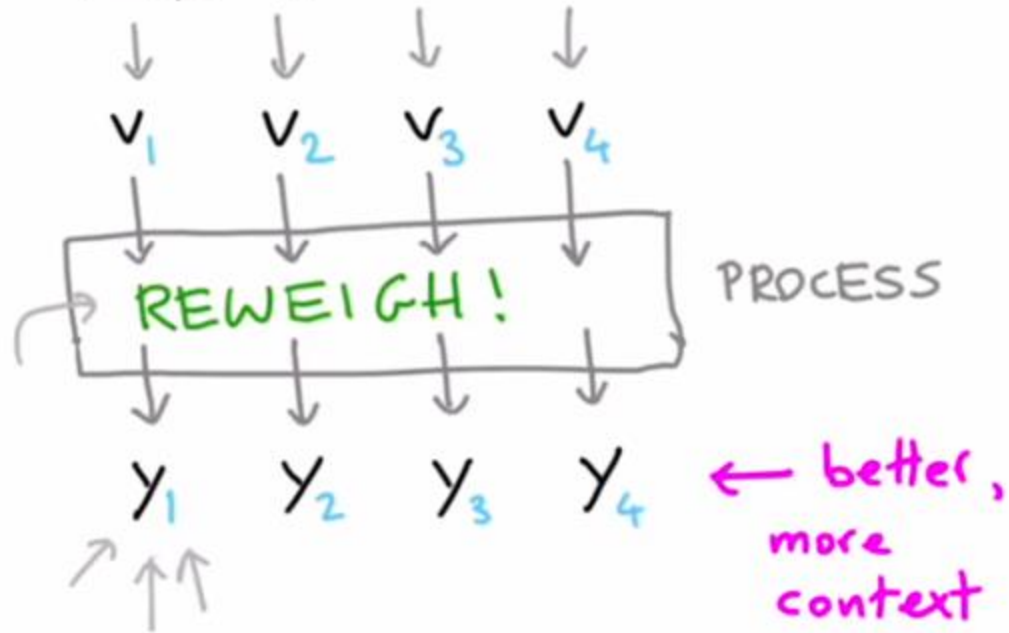https://learning.rasa.com/transformers/self-attention/

Bank of the river.

$\downarrow$ $\downarrow$ $\downarrow$ $\downarrow$

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

REWEIGH!    PROCESS

$y_1$ $y_2$ $y_3$ $y_4$ $\leftarrow$ better, more context

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

$V_1 V_1 = W_{11}$            $W_{11}$

$V_1 V_2 = W_{12}$ $\xrightarrow{\hspace{2cm}}$ $W_{12}$

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

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

$\longrightarrow$ sum to one

Bank of the river.

$\downarrow$ $\downarrow$ $\downarrow$ $\downarrow$

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

REWEIGH!   PROCESS

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

$V_1 V_1 = W_{11}$

$V_1 V_2 = W_{12}$ ⟶ norm

$V_1 V_3 = W_{13}$

$V_1 V_4 = W_{14}$

⌐ sum to one

$W_{11}$

$W_{12}$

$W_{13}$

$W_{14}$

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

↑       ↑       ↑       ↑

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

$V_1 V_1 = W_{11}$  $\qquad W_{11}$
$V_1 V_2 = W_{12}$  $\xrightarrow{\quad}$  $W_{12}$
$V_1 V_3 = W_{13}$  $\quad$ norm  $\quad W_{13}$
$V_1 V_4 = W_{14}$  $\qquad W_{14}$

↳ 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.

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

REWEIGH!  PROCESS

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

$V_1 V_1 = W_{11}$     $W_{11}$

$V_1 V_2 = W_{12}$  ⟶  $W_{12}$

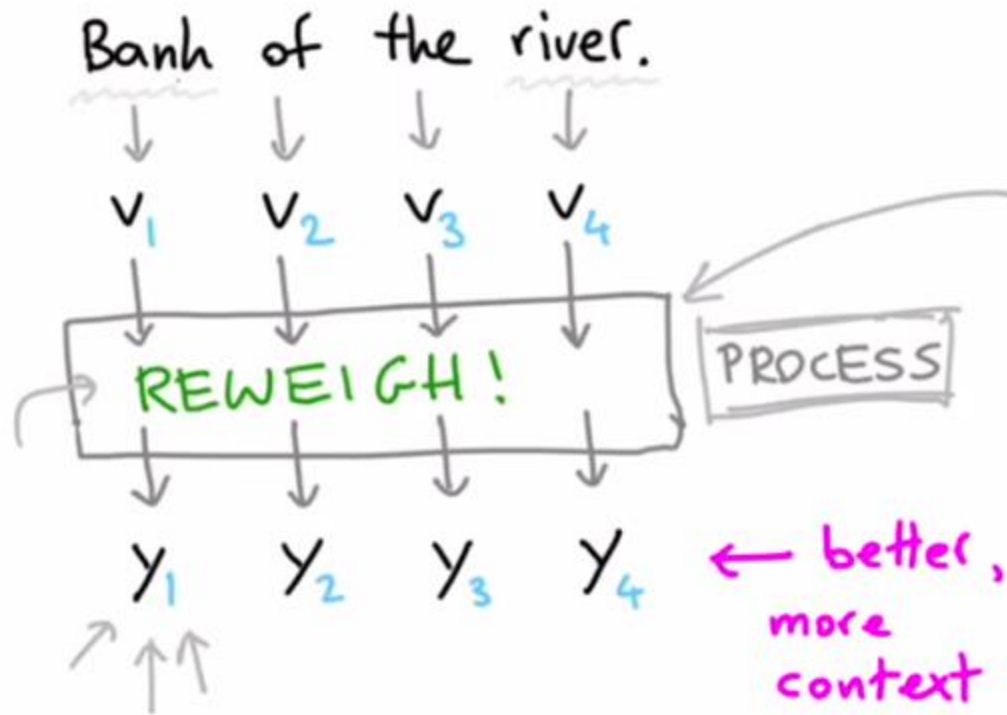$V_1 V_3 = W_{13}$    norm   $W_{13}$

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

↑ sum to one

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

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

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

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

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

Bank 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
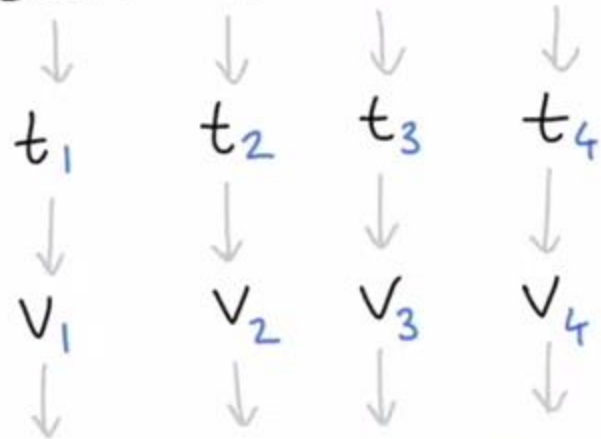- Proximity has no influence
- Shape independant

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

norm →

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

sum to one

SELF ATTENTION

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

$V_1$    $V_2$    $V_3$    $V_4$      $Y_4$

$Y_3$

$Y_2$

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

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

$Y_1$   $Y_2$   $Y_3$   $Y_4$

$S_{31}$    $S_{32}$    $S_{33}$    $S_{34}$     no weights $\longrightarrow$ **ATTENTION**

**DOT PRODUCT**    $\leftarrow V_3$      $V_1$   $V_2$   $V_3$   $V_4$

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

26

$V_1$   $V_2$   $V_3$   $V_4$   $\dashrightarrow$ $Y_4$

$Y_3$

$Y_2$

$Y_1$

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

$\sum_j W_{3j} = 1$

NORMALISE

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

~~no~~ weights $\longrightarrow$

$Y_1$   $Y_2$   $Y_3$   $Y_4$

ATTENTION

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

DOT PRODUCT   $\leftarrow V_3$

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

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

$\longrightarrow Y_3$

$\longrightarrow Y_2$

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

NORMALISE

$\sum_j W_{3j} = 1$

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

DOT PRODUCT $\quad \leftarrow V_3$

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

30

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

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

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

NORMALISE

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

DOT PRODUCT $\leftarrow V_3$

QUERY

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

$Y_4$

$Y_3$

$Y_2$

$Y_1$

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

$\rightarrow y_3$

$\rightarrow y_2$

$W_{31} \quad W_{32} \quad W_{33} \quad W_{34} \qquad \rightarrow y_1$

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

NORMALISE

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

QUERY

DOT PRODUCT $\quad \leftarrow V_{\underline{3}}$

$V_1 \quad V_2 \quad V_3 \quad V_4 \qquad \leftarrow$ KEYS

$V_1$    $V_2$    $V_3$    $V_4$      $Y_4$

$Y_3$

$Y_2$

$Y_1$

VALUES

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

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

NORMALISE

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

QUERY

DOT PRODUCT    $\leftarrow V_3$

KEYS

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

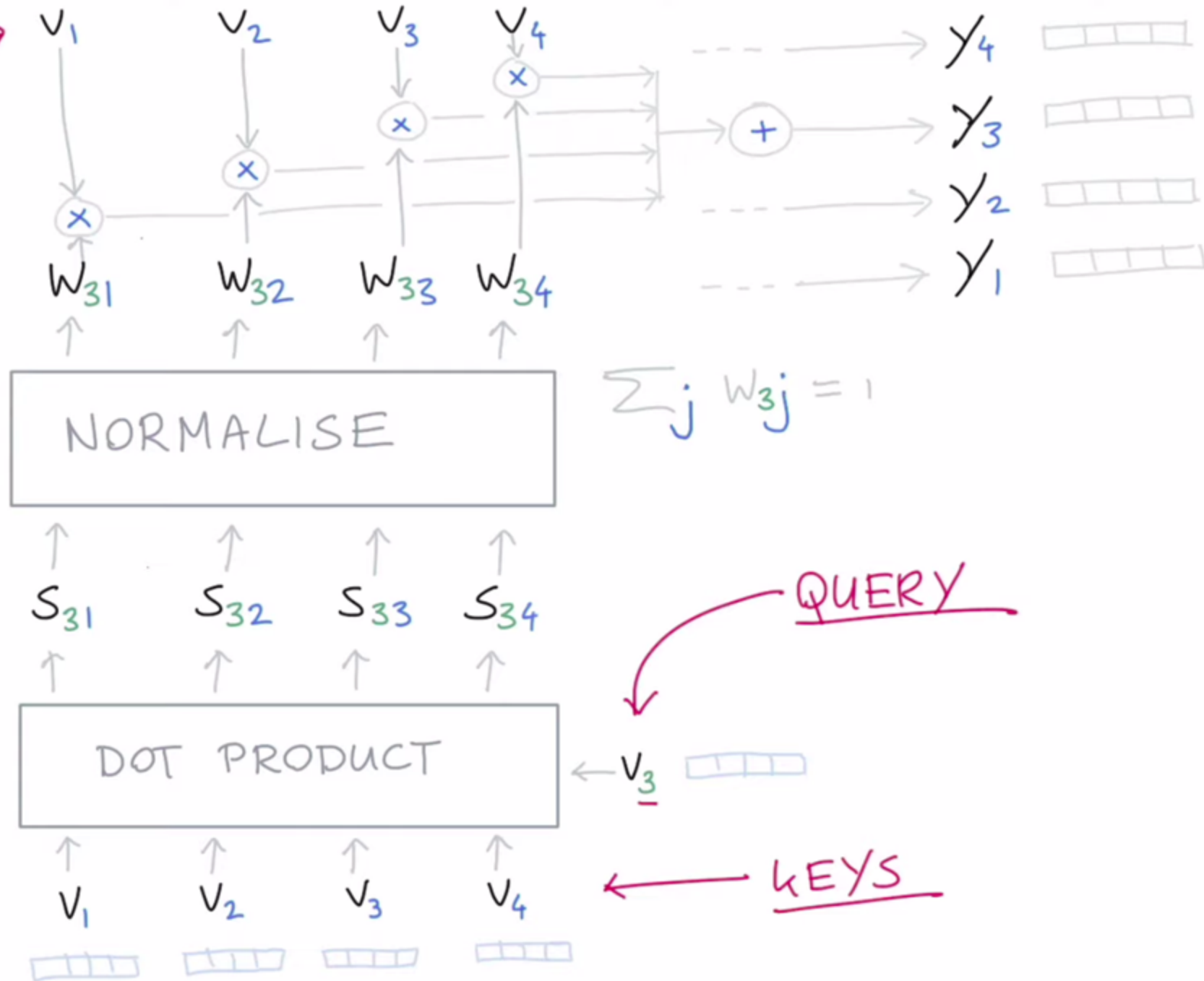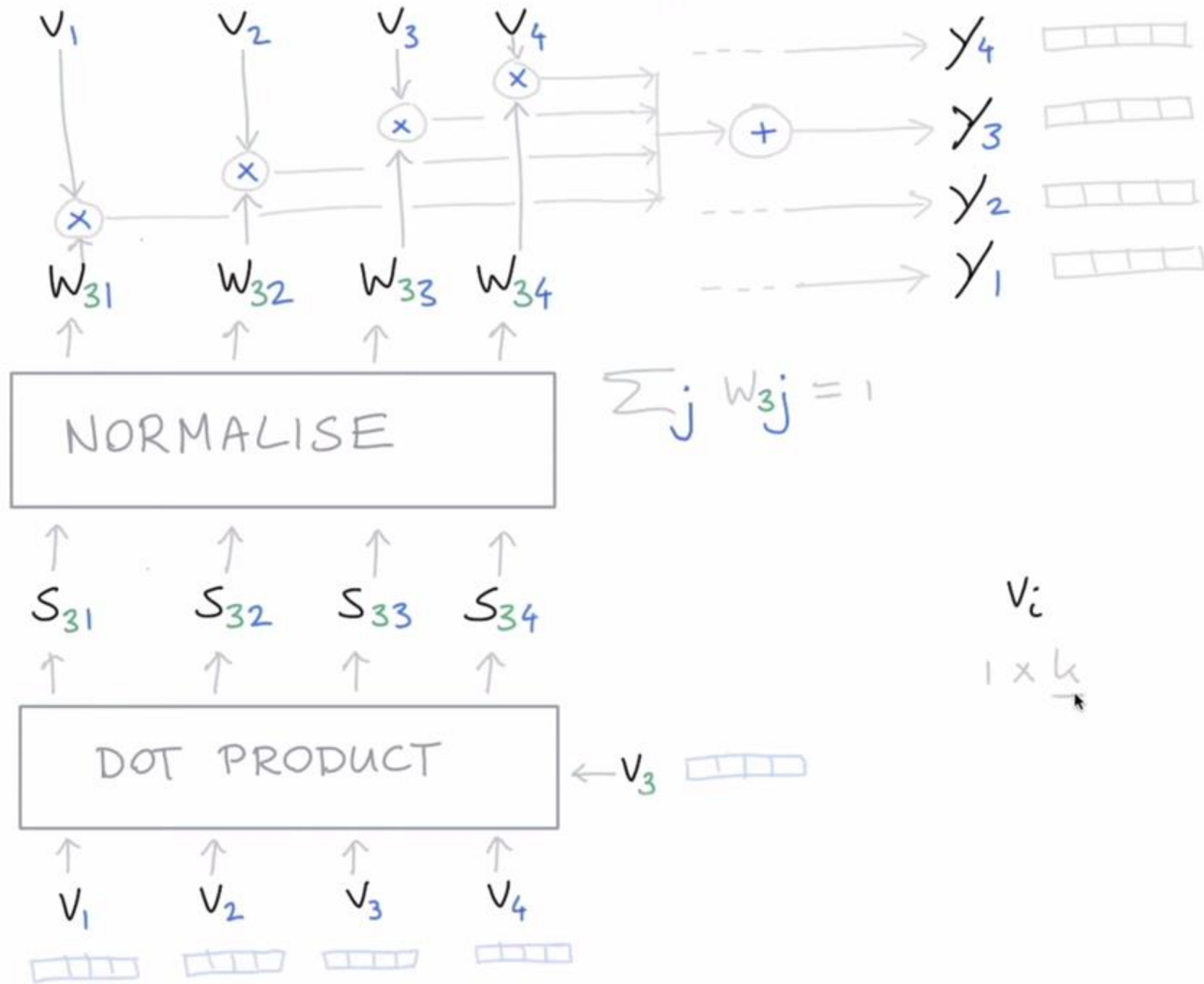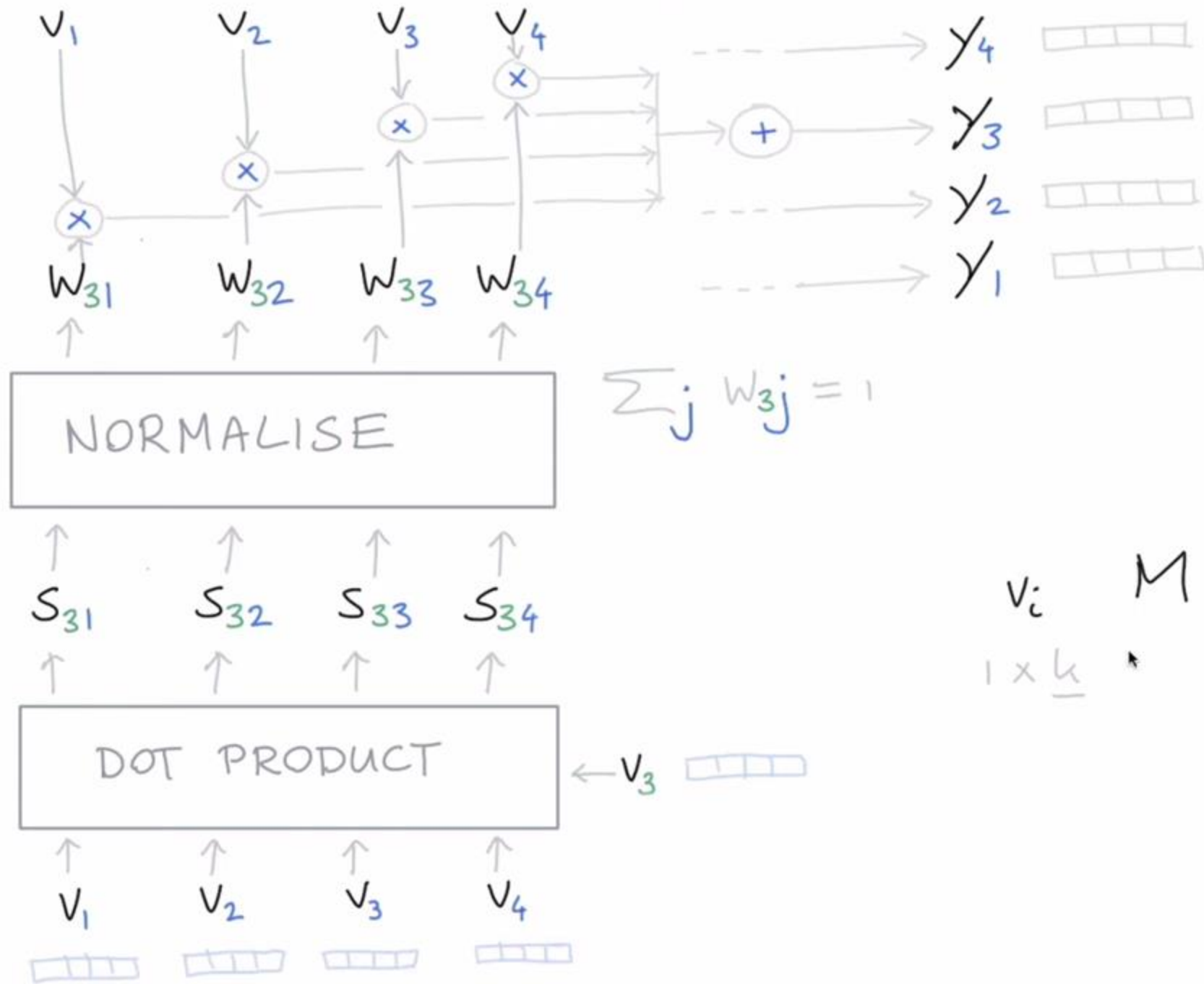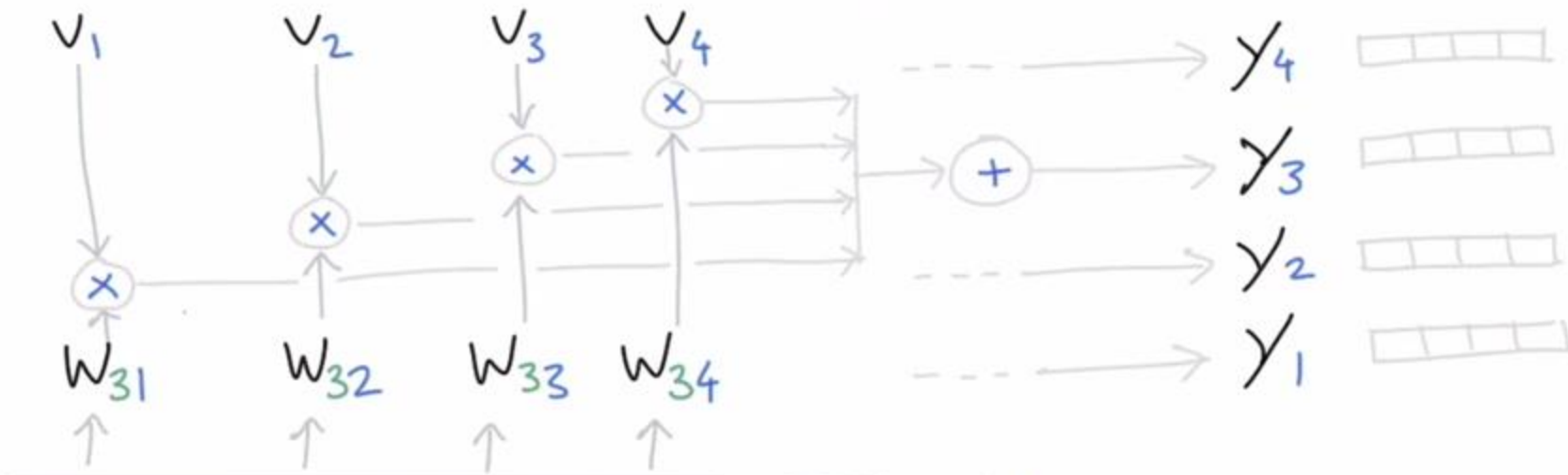# Key, Value and Query

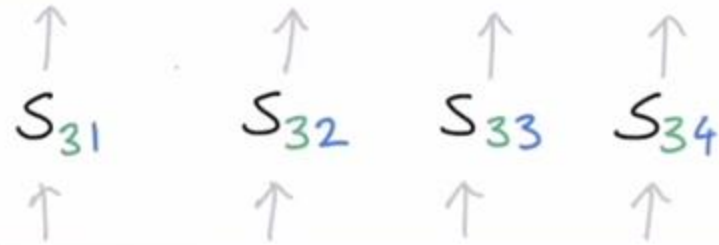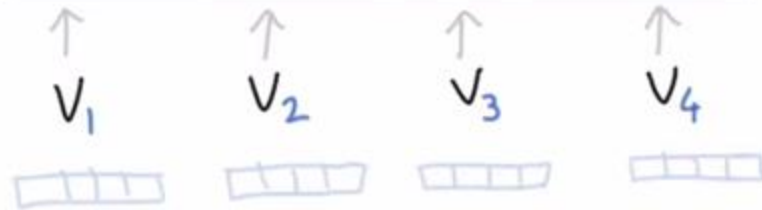- The major component in the transformer is the unit of <span style="color:red">multi-head self-attention mechanism.</span>

$V_1$   $V_2$   $V_3$   $V_4$   → $Y_4$

→ $Y_3$

$W_{31}$   $W_{32}$   $W_{33}$   $W_{34}$   → $Y_2$

→ $Y_1$

NORMALISE

$\sum_j W_{3j} = 1$

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

$V_i$   $M$

$1 \times \underline{k}$

DOT PRODUCT   ← $V_3$

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

36

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

$\times$  $\times$  $\times$  $\times$

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

NORMALISE

$\sum_j W_{3j} = 1$

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

DOT PRODUCT  $\leftarrow V_3$

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

$Y_4$
$Y_3$
$Y_2$
$Y_1$

$V_i$  $M$

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

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

$\to y_3$

$\to y_2$

$\to y_1$

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

$\sum_j W_{3j} = 1$

NORMALISE

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

$v_i \quad M = [\ ]$

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

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

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

DOT PRODUCT  $\leftarrow V_3$

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

$V_i$  $M = []$

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

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

$y_3$

$y_2$

$y_1$

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

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

NORMALISE

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

$v_i \qquad M = [\ ]$

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

DOT PRODUCT

$\leftarrow v_3$

$v_1 M_k \quad v_2 M_k \quad v_3 M_k \quad v_4 M_k \quad kEY$

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

$y_3$

$y_2$

$y_1$

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

$\sum_j W_{3j} = 1$

NORMALISE

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

QUERY

$v_i \qquad M = [\ ]$

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

DOT PRODUCT

$\leftarrow v_3 M_Q$

$v_1 M_k \quad v_2 M_k \quad v_3 M_k \quad v_4 M_k \quad$ KEY

42

VALUE

$V_1 M_v$  $V_2 M_v$  $V_3 M_v$  $V_4 M_v$  — — — → $Y_4$  ☐☐☐☐

⊗ → $Y_3$

⊗

⊗

⊗  (+) → $Y_3$  ☐☐☐☐

— — — → $Y_2$  ☐☐☐☐

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

↑  ↑  ↑  ↑

$\sum_j W_{3j} = 1$

NORMALISE

↑  ↑  ↑  ↑

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

↑  ↑  ↑  ↑

QUERY

$V_i$  $M$  $=$  $[\,]$

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

DOT PRODUCT  ← $V_3 M_Q$ ☐☐☐

↑  ↑  ↑  ↑

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

☐☐☐  ☐☐☐  ☐☐☐  ☐☐☐

43

SELF ATTENTION

KEYS
LINEAR

QUERIES
LINEAR

VALUES
LINEAR

MATMUL → $[S_{ij}]$ ↑ NORMALISE → $[W_{ij}]$

MATMUL

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

$v_i$ ▭▭▭

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

44

45

SELF ATTENTION

KEYS
LINEAR

QUERIES
LINEAR

VALUES
LINEAR

MATMUL

$[S_{ij}]$

NORMALISE

$\rightarrow [W_{ij}]$

MATMUL

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

$v_i \; \square\square\square$

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

NER

SELF ATTENTION

SELF ATTENTION

$v_1 \qquad v_2 \qquad v_3$

SELF ATTENTION

KEYS

LINEAR

QUERIES

LINEAR

VALUES

LINEAR

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

$v_i$

MATMUL

$[S_{ij}]$

NORMALISE

$\rightarrow [W_{ij}]$

MATMUL

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

NER

SELF ATTENTION

SELF ATTENTION

$V_1$   $V_2$   $V_3$

SELF ATTENTION

KEYS

LINEAR

QUERIES

LINEAR

VALUES

LINEAR

MATMUL

NORMALISE

$[S_{ij}]$

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

$[W_{ij}]$

$v_i$

MATMUL

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

NER

SELF ATTENTION

SELF ATTENTION

$v_1$   $v_2$   $v_3$

SELF ATTENTION

KEYS
LINEAR

QUERIES
LINEAR

VALUES
LINEAR

MATMUL

NORMALISE

$[S_{ij}]$

$[W_{ij}]$

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

$V_1 \quad V_2 \quad V_3$

49

SELF ATTENTION

KEYS
LINEAR

QUERIES
LINEAR

VALUES
LINEAR

MATMUL → NORMALISE

$[S_{ij}]$

$\to [W_{ij}]$

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

$v_1$   $v_2$   $v_3$

50

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

$V_1$ $V_2$ $V_3$

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

SELF ATTENTION

KEYS
LINEAR

QUERIES
LINEAR

VALUES
LINEAR

$[S_{ij}]$

MATMUL

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

$V_1$    $V_2$    $V_3$

52

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
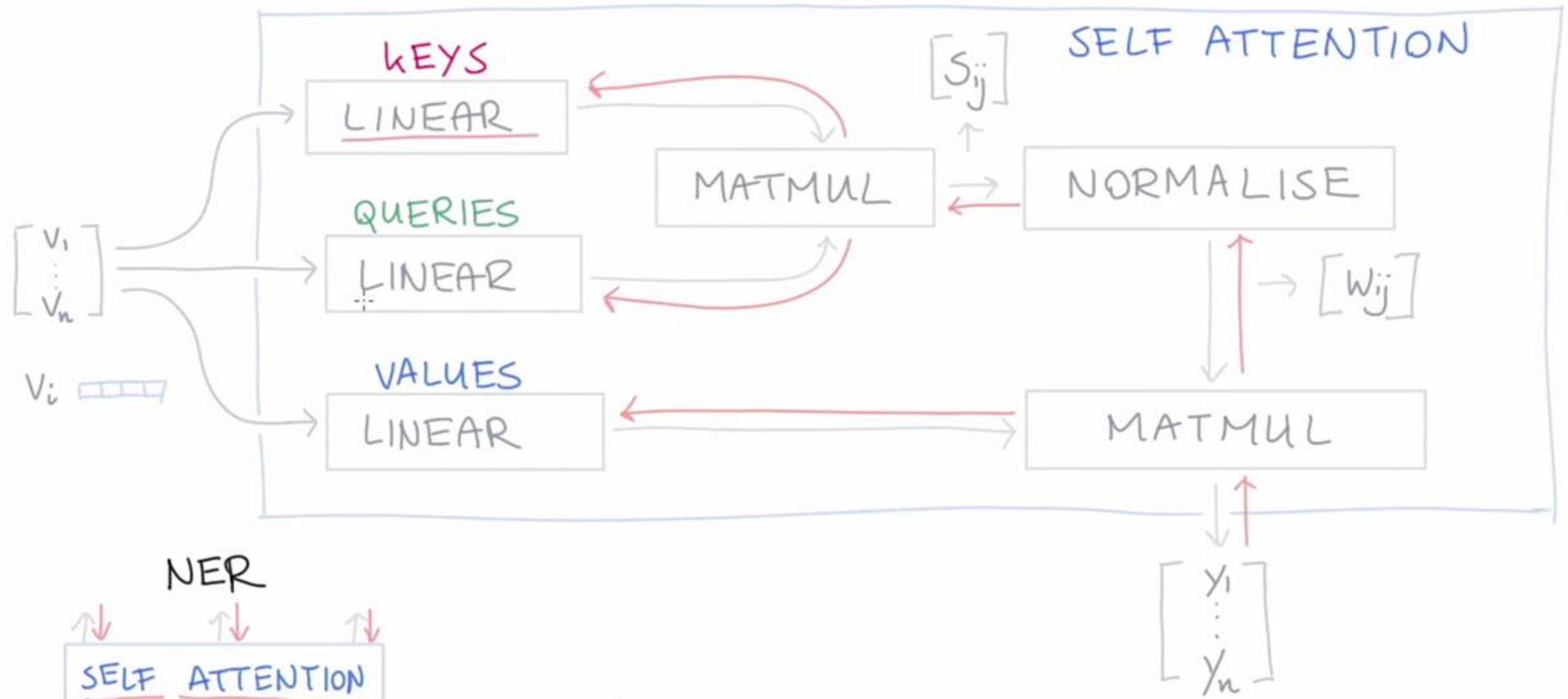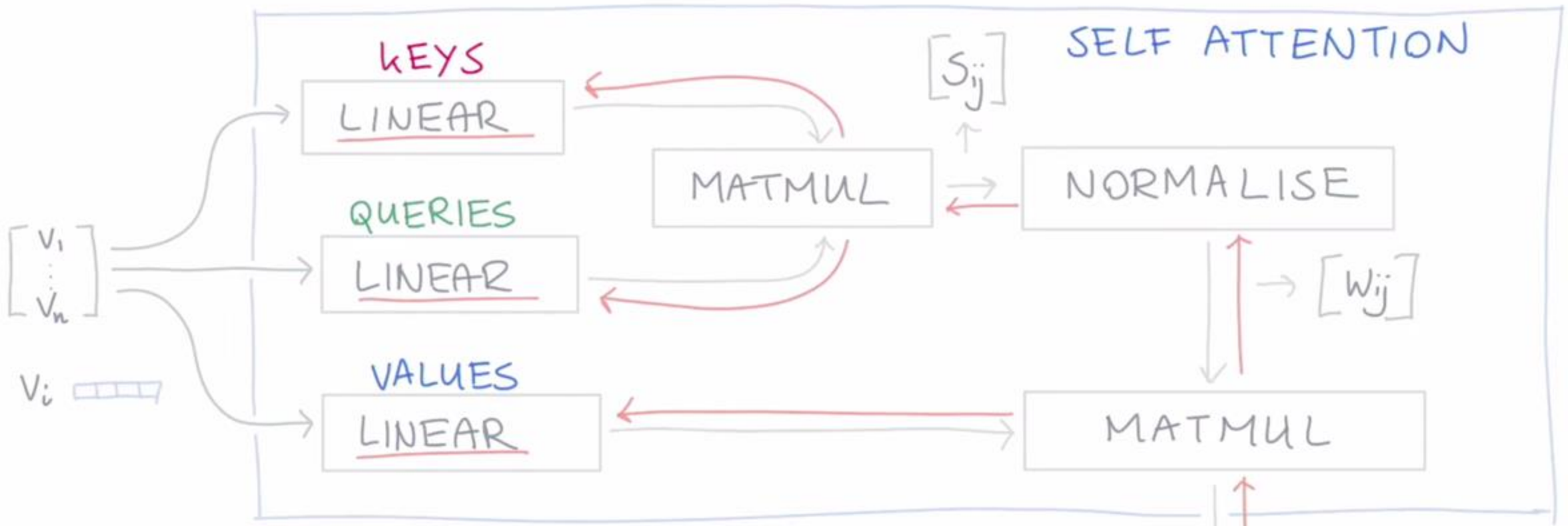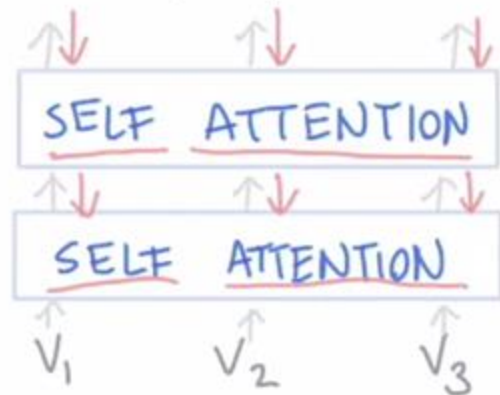
SELF ATTENTION

$V_1$    $V_2$    $V_3$

← clich together ⚓

53

# 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

**Łukasz Kaiser**[*]
Google Brain
lukaszkaiser@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.

[*]Equal contribution. Listing order is random. Jakob proposed replacing RNNs with self-attention and started the effort to evaluate this idea. Ashish, with Illia, 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

# The annotated Transformer (Code)

- [The Annotated Transformer (harvard.edu)](harvard.edu)

# Acknowledgments

These slides were adapted from the book

# 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/

- 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