Lecture Notes for

# Neural Networks and Machine Learning

Transformers and Vision Transformers

# Logistics and Agenda

- Logistics
  - Paper presentations (Thurs)
- Agenda
  - Transformers
- Next Time:
  - Vision Transformers
  - Paper Presentation
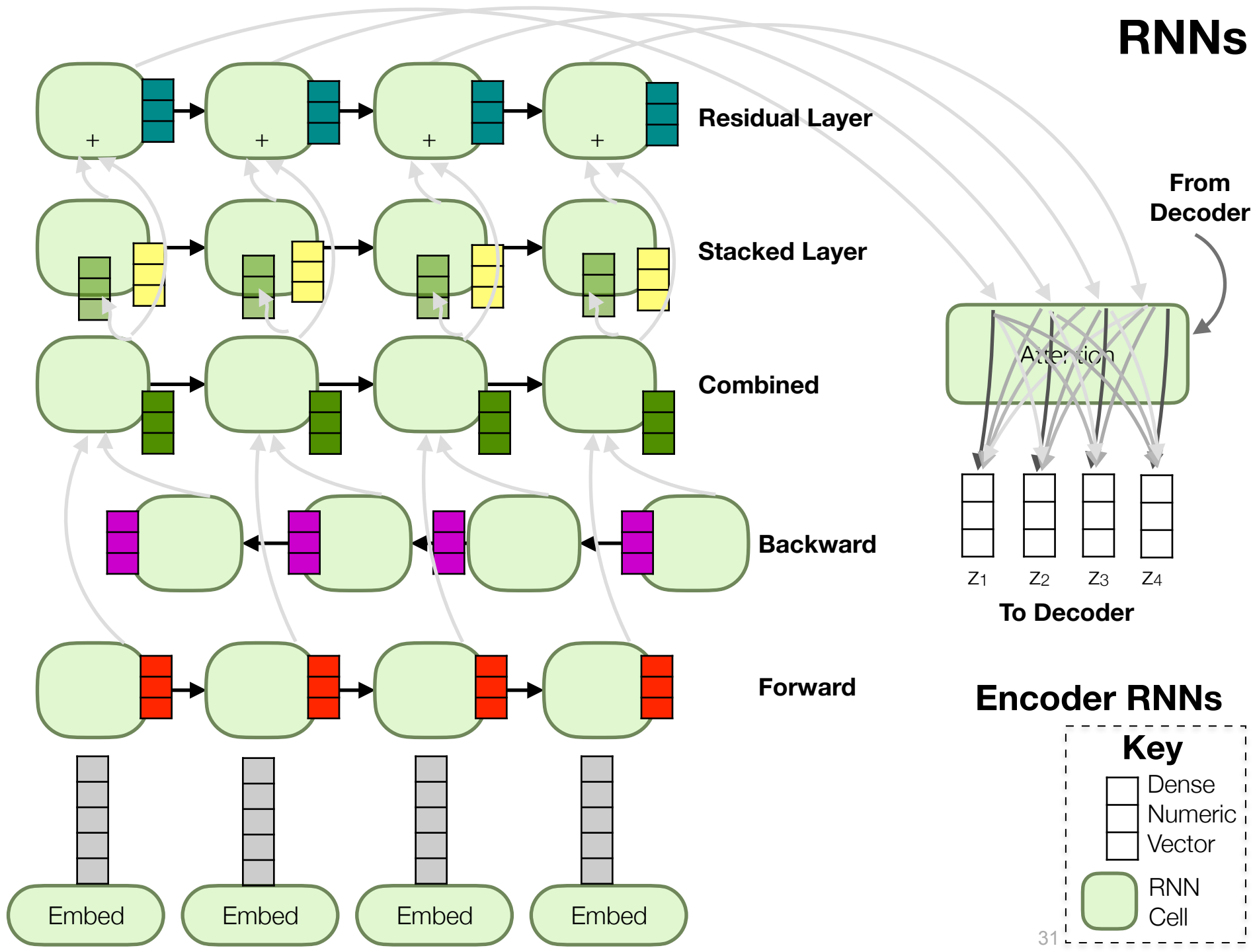  - Consistency losses

# Transformers

**Dr Simone Stumpf** @DrSimoneS... · 13h  ···
God grant me the confidence of an average
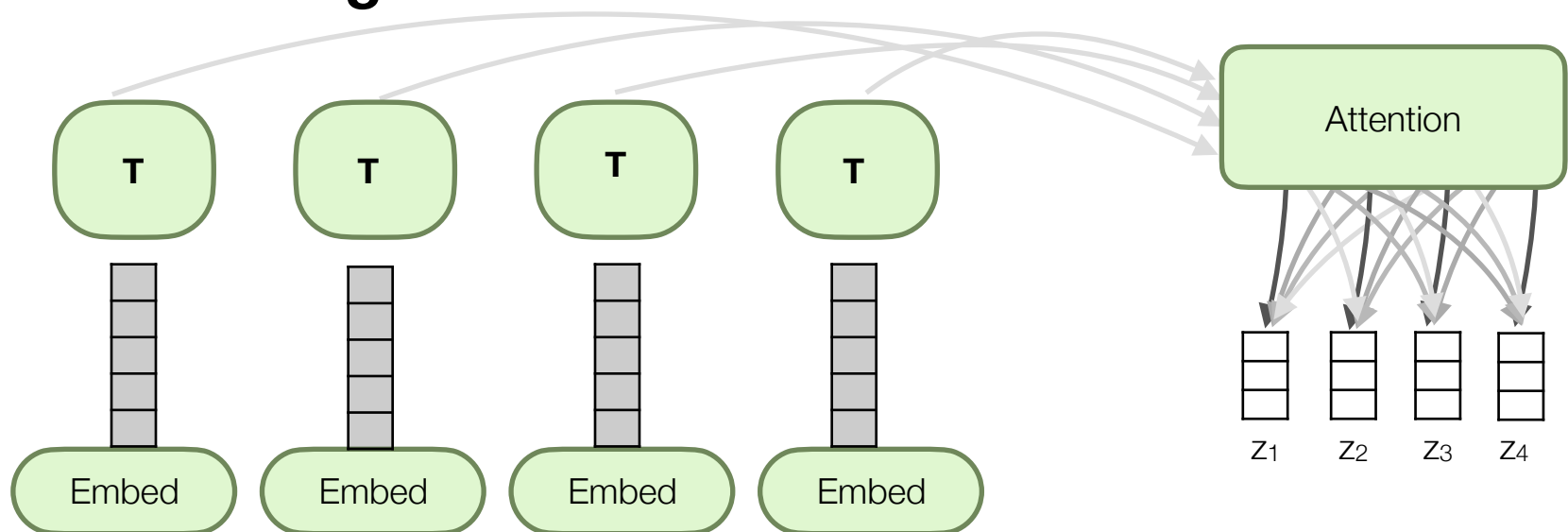machine learning expert.

# RNNs

**Residual Layer**

**Stacked Layer**

**Combined**

**From Decoder**

Attention

**Backward**

$Z_1$    $Z_2$    $Z_3$    $Z_4$

**To Decoder**

**Forward**

## Encoder RNNs

Embed     Embed     Embed     Embed

**Key**

Dense
Numeric
Vector

RNN
Cell

31

# Transformers Intuition

- Recurrent networks track state using an "updatable" state vector, but this takes lots of processing to across sequence

- Attention mechanism (in RNNs) already takes a weighted sum of state vectors to generate new token in a decoder

- … so why not just use attention on a transformation of the embedding vectors? **Do away with the recurrent state vector all together?**

# Attention is All You Need
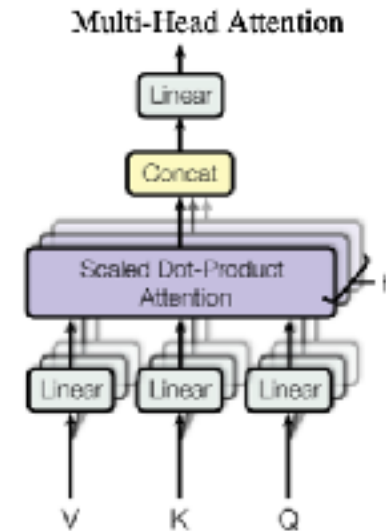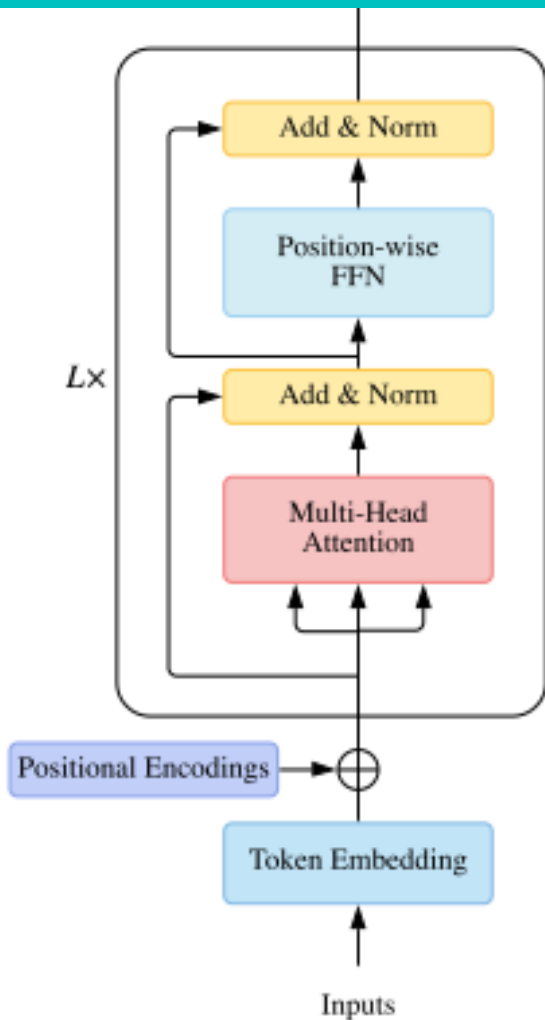
- **Continued Motivation**:
  - RNNs are not inherently parallelized or efficient at remembering based on state vector
  - CNNs are not resilient to long-term word relationships, limited by filter size
- **Transformer Solution:**
  - Build attention into model from the **beginning**
  - Compare all words to each other through **self-headed** attention
  - Define a notion of "**position**"in the sequence
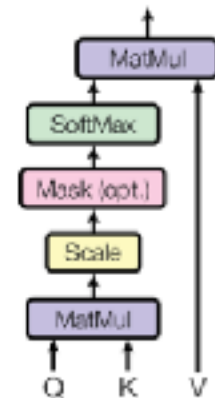  - *Should be resilient to long term relationships and be highly parallelized for GPU computing!!*

# Transformer Overview



Multi-Head Attention

more than one
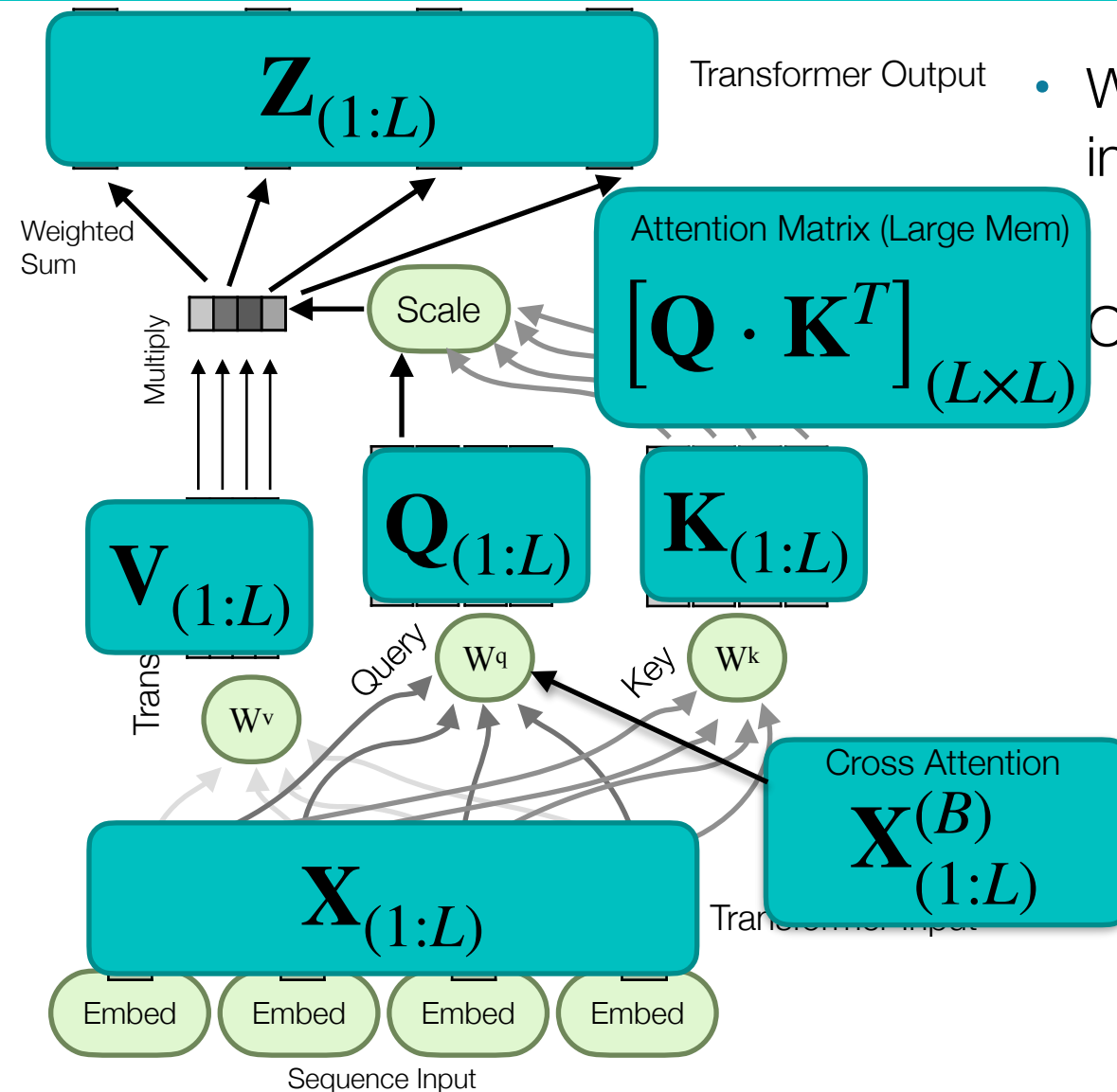Q,K,V use in document

Scaled Dot-Product Attention

for each word

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$

$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

https://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf

# Self Attention Overview



**Z**$_{(1:L)}$ — Transformer Output

Weighted Sum

Multiply

Scale

Attention Matrix (Large Mem)

$$\left[ \mathbf{Q} \cdot \mathbf{K}^T \right]_{(L \times L)}$$

**V**$_{(1:L)}$

**Q**$_{(1:L)}$

**K**$_{(1:L)}$

Trans

W$^v$

W$^q$

W$^k$

Query

Key

Cross Attention

**X**$^{(B)}_{(1:L)}$

**X**$_{(1:L)}$

Transformer Input

Embed   Embed   Embed   Embed

Sequence Input

- What parameters are trained in diagram?

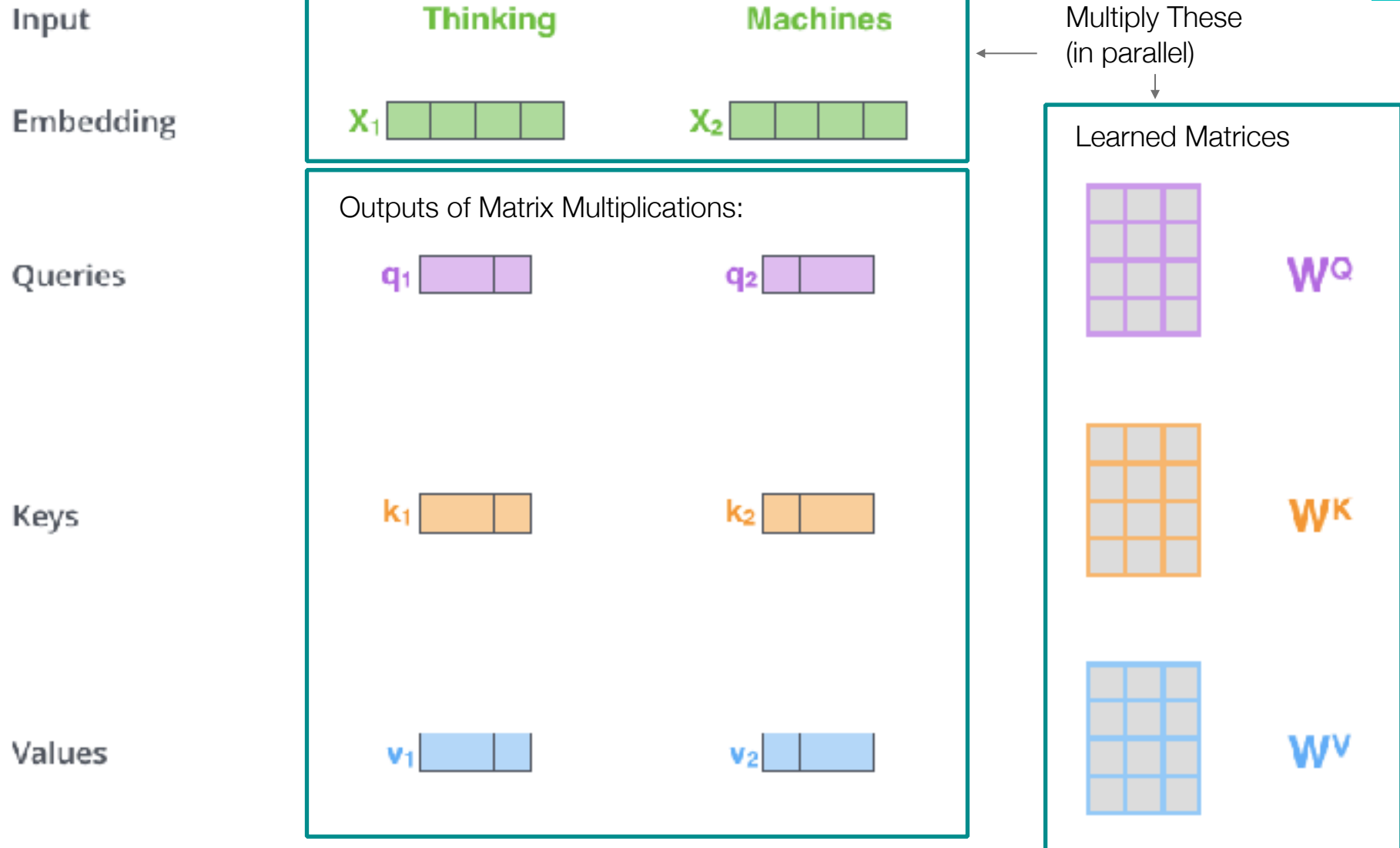- $\mathbf{W}^v, \mathbf{W}^q, \mathbf{W}^k$

Other Parameters:

- $L$: length of sequence

- Query/Key dimension, $d_k$

- Value dimension, $d_v$

- How many times to apply attention (i.e., number of heads)

- Type of positional encoding (more later)

# Transformer: in more detail

| | Thinking | Machines | |
|---|---|---|---|
| Input | | | Multiply These (in parallel) |

Embedding: $X_1$ [ ][ ][ ][ ]   $X_2$ [ ][ ][ ][ ]

Outputs of Matrix Multiplications:

Queries: $q_1$ [ ][ ]   $q_2$ [ ][ ]   $W^Q$

Keys: $k_1$ [ ][ ]   $k_2$ [ ][ ]   $W^K$

Values: $v_1$ [ ][ ]   $v_2$ [ ][ ]   $W^V$

Learned Matrices

**Excellent Blog on Transformers:** http://jalammar.github.io/illustrated-transformer/

# Transformer: in more detail



Input
Embedding
Queries
Keys
Values
Score
Divide by 8 ( $\sqrt{d_k}$ )
in visual, $d_k = 3$
Softmax
Softmax
  X
Value
Sum

**Thinking**          **Machines**

$x_1$          $x_2$

Calc. q, k, v for each word

$q_1$          $q_2$
$k_1$          $k_2$
$v_1$          $v_2$

$q_1 \cdot k_1 = 112$          $q_1 \cdot k_2 = 96$

Divide          Divide
14          12

Softmax

0.88          0.12

Multiply          Calc weights          Multiply

$v_1$          $v_2$

weighted sum for all words in document

Sum

$z_1$          $z_2$
attention for word 1          attention for word 2

Straight forward to do this operation in matrix form:

X          $W^Q$          Q
Thinking
Machines          $\leftarrow d_k \rightarrow$          $q_1$ $q_2$

X          $W^K$          K
Thinking
Machines          $\leftarrow d_k \rightarrow$          $k_1$ $k_2$

X          $W^V$          V
Thinking
Machines          $\leftarrow d_v \rightarrow$          $v_1$ $v_2$

Q          $K^T$          V

$$\text{softmax}\left( \frac{\phantom{xx} \times \phantom{xx}}{\sqrt{d_k}} \right)$$

Z
$=$          $z_1$ $z_2$

Size of W matrices:
$W^V$: |Embed Size| x $d_v$
$W^{Q,K}$: |Embed Size| x $d_k$

Size of Q,K,V:
|Seq Len| x $d_v$ or $d_k$

**Excellent Blog on Transformers:** http://jalammar.github.io/illustrated-transformer/

# Transformer: Multi-headed Attention



one row
for each
word

Num columns
determined by
$W^O$

$(words \times Z_{concat})$  $\times W^O = Z$

$(Z_{concat} \times Embed)$    $(words \times Embed)$

**Excellent Blog on Transformers:** http://jalammar.github.io/illustrated-transformer/

# Putting It Together



Prediction

Fully Connected

Global Average (or other)

Layer Normalization

Fully Connected

Layer Normalization
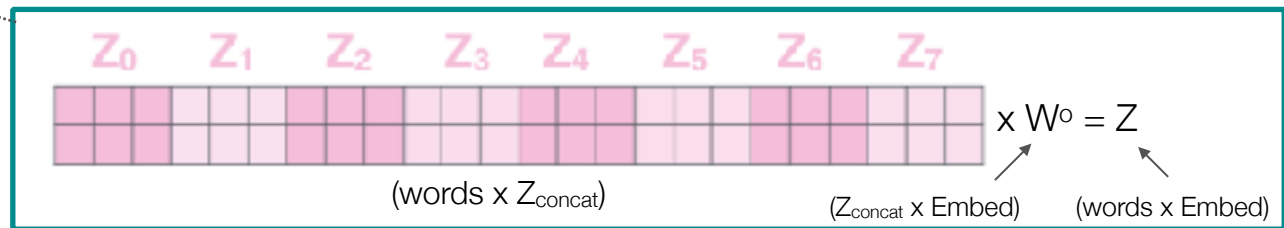
```
tf.keras.layers.MultiHeadAttention(
    num_heads,       (Number of heads $Z_1$-$Z_7$)
    key_dim,         (size of query/key $d_k$)
    value_dim,       (size of each $d_v$)
    output_shape,    (Embed size of Z)
    …
```
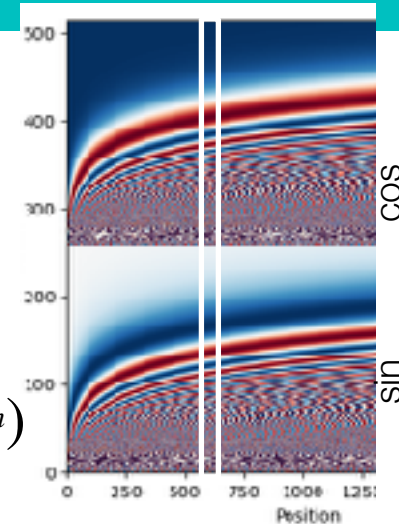
$$LN(Z_{row}) = \alpha \frac{Z_i - \mu_Z}{\sqrt{\sigma_Z^2 + \epsilon}} + \beta$$

Learn to normalize the rows of Z

$Z_0$  $Z_1$  $Z_2$  $Z_3$  $Z_4$  $Z_5$  $Z_6$  $Z_7$

x W$^o$ = Z

(words x $Z_{concat}$)         ($Z_{concat}$ x Embed)    (words x Embed)

?

# Transformer: Positional Encoding

- Objective: add notion of position to embedding
- Attempt in paper: add sin/cos to embedding

$p$: in sequence
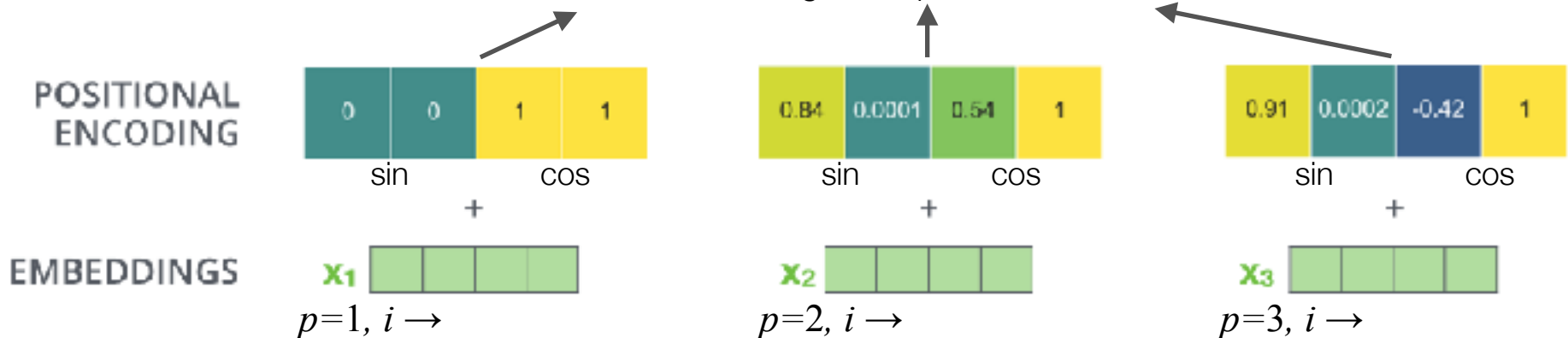$d\_m$: 1/2 dim of embed
$i$ = index in vector

$$PE_{(p,i\in 0...d_m-1)} = \sin(p/10000^{i/d_m})$$

$$PE_{(p,i\in d_m...2d_m)} = \cos(p/10000^{(i-d_m)/d_m})$$



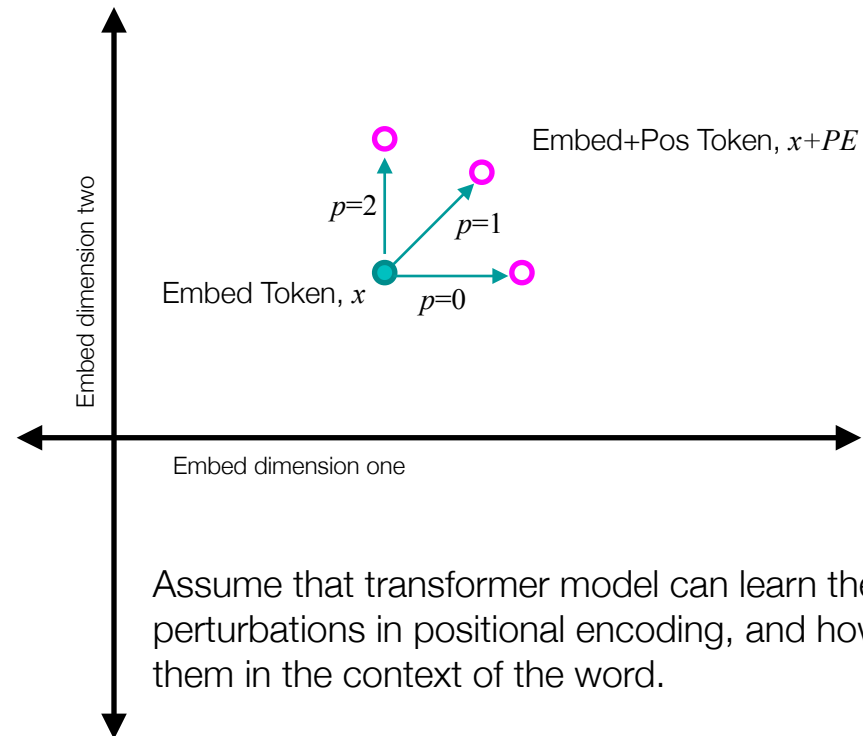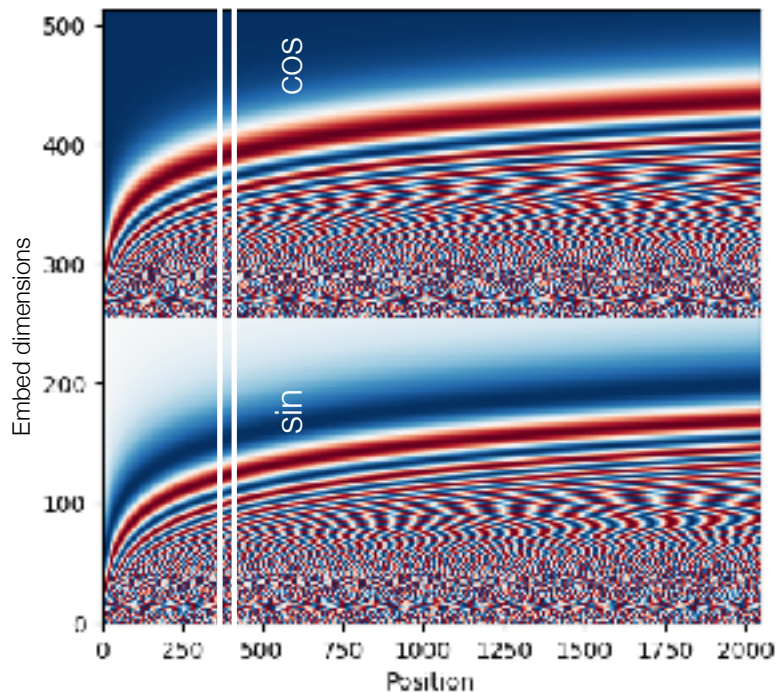Now use the new embeddings, with position, into transformer architecture

POSITIONAL
ENCODING

| 0 | 0 | 1 | 1 |

   sin       cos

+

| 0.84 | 0.0001 | 0.54 | 1 |

   sin       cos

+

| 0.91 | 0.0002 | -0.42 | 1 |

   sin       cos

+

EMBEDDINGS

$X_1$    $X_2$    $X_3$

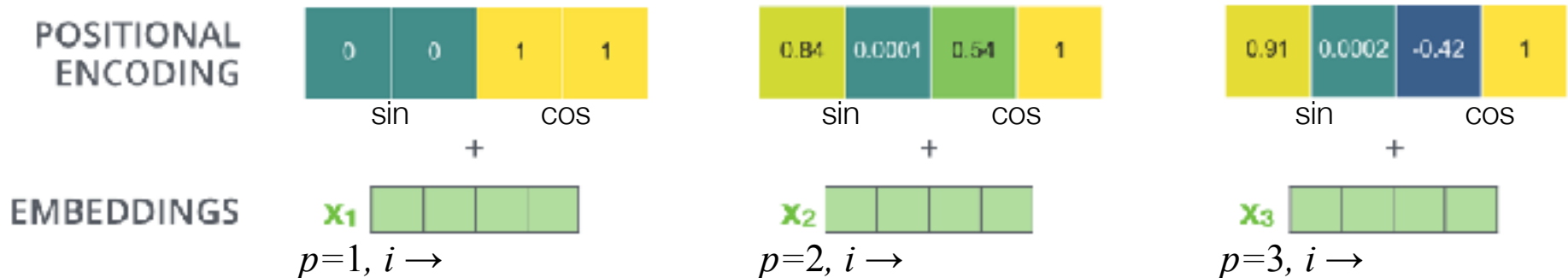$p=1, i \rightarrow$      $p=2, i \rightarrow$      $p=3, i \rightarrow$

**Hypothesis**: Now the word proximity is encoded in the embedding matrix, with other pertinent information.  Well, it does help… so it could be true that this is a good way to do it.

**Excellent Blog on Transformers:** http://jalammar.github.io/illustrated-transformer/
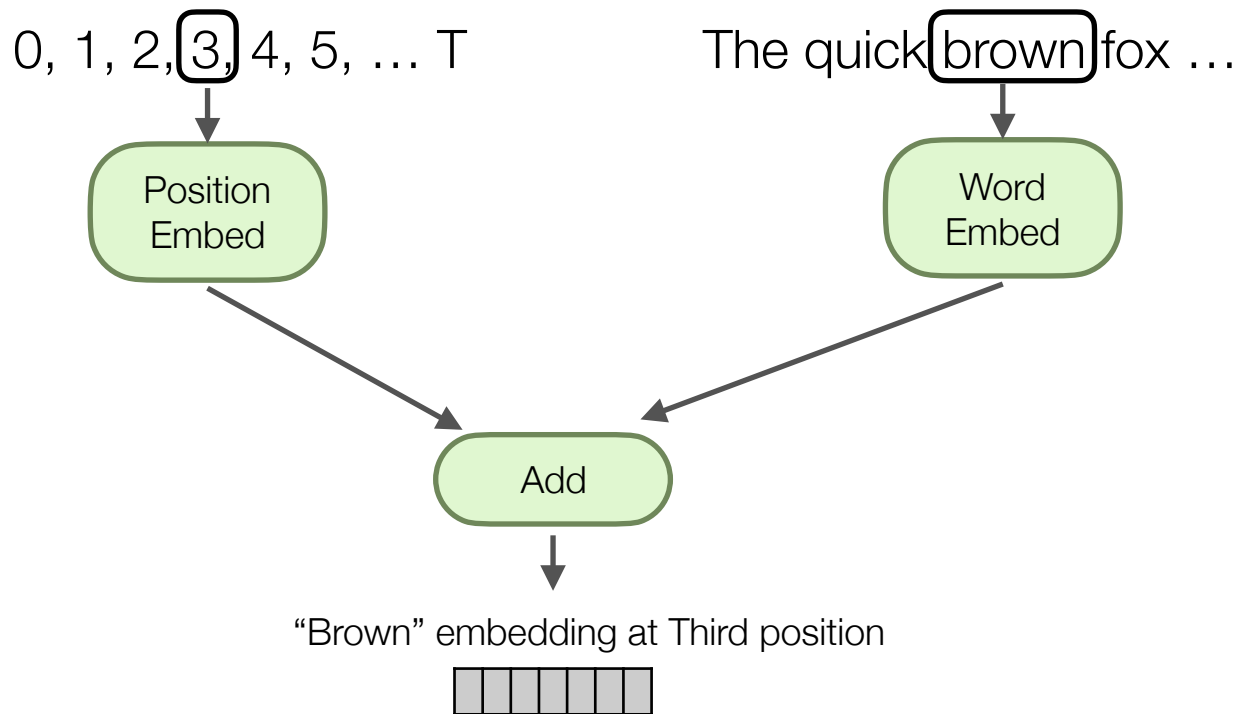
# Positional Intuition, Geometrically



Assume that transformer model can learn the small perturbations in positional encoding, and how to use them in the context of the word.

Embed+Pos Token, $x+PE$

Embed Token, $x$

$p=2$ $p=1$ $p=0$

Embed dimension two

Embed dimension one

POSITIONAL ENCODING

| 0 | 0 | 1 | 1 |
|---|---|---|---|
| sin | | cos | |

+

EMBEDDINGS

$X_1$

$p=1, i \rightarrow$

| 0.84 | 0.0001 | 0.54 | 1 |
|---|---|---|---|
| sin | | cos | |

+

$X_2$

$p=2, i \rightarrow$

| 0.91 | 0.0002 | -0.42 | 1 |
|---|---|---|---|
| sin | | cos | |

+

$X_3$

$p=3, i \rightarrow$

# Transformer: Positional Encoding

- Objective: add notion of position to embedding
- Attempt in original paper: add sin/cos to embedding
- **But could be anything that encodes position, like:**

0, 1, 2, 3, 4, 5, … T          The quick brown fox …

Position Embed

Word Embed

Add

"Brown" embedding at Third position

Lecture Notes for

# Neural Networks and Machine Learning

Transformers

**Next Time:**
SSL, Vision Transformers

**Reading:** None