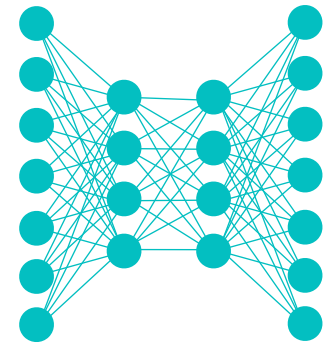


Lecture Notes for **Neural Networks and Machine Learning**



Transformers



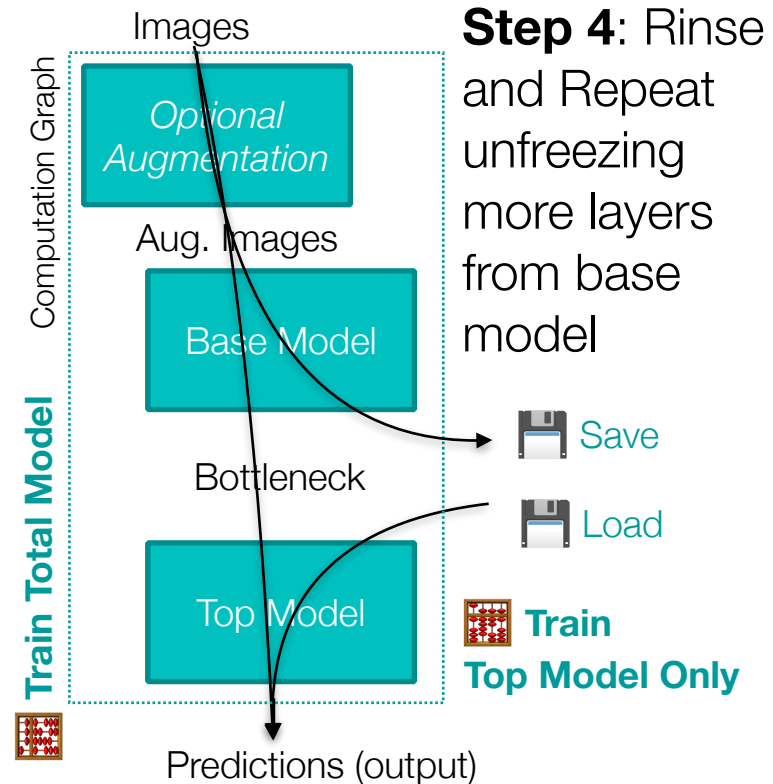
Logistics and Agenda

- Logistics
 - Paper presentations
- Agenda
 - Finish Transfer Demo
 - Transformers
- Next Time:
 - Position Encoding Transformers
 - Vision Transformers
 - Paper Presentation
 - Consistency losses



Freezing and Fine-tuning Efficiently

- **Step 1:** Freeze entire base model:
 - No update during back-propagation
 - *Optional: Augment a set of training data*
 - Send training dataset through base model
 - ◆ Save out bottleneck features
- **Step 2:** Train bottleneck features in new task
 - Typically 5-10 epochs is sufficient, easy to overfit (very fast)
 - Larger training step size is okay
- **Step 3:** Fine-tune, **unfreeze** a few layers in base model:
 - Attach newly trained model to pre-trained model, *Optional: use augmentation*
 - Train to your hearts content, use smaller training step size



```
images = Input(shape=(IMG_SIZE, IMG_SIZE, 3))
base_model = VGG(include_top=False,
                 input_tensor=images,
                 weights="imagenet")

bottleneck = Input(shape=base_model.output.shape)
predict = Dense(NUM_CLASSES)(bottleneck)
top_model = Model(bottleneck, predict)

model_total = Model(images, outputs=
top_model(base_model.output))
```





Bottlenecking on a GPU

Dogs versus Cats



justinledford Justin Ledford

 Member of [8000net](#)



eclarson Eric Larson

Updated for `tf==2.12` in the Main Repository:
[02 Transfer Learning.ipynb](#)

Original Example: <https://github.com/8000net/Transfer-Learning-Dolphins-and-Sharks>

Another Great Example:

https://keras.io/examples/vision/image_classification_efficientnet_fine_tuning/



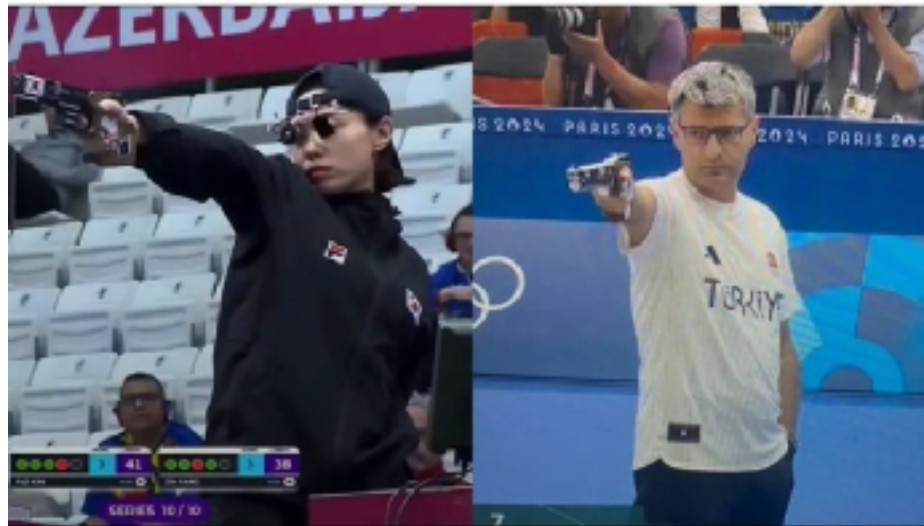
Justin Ledford ·



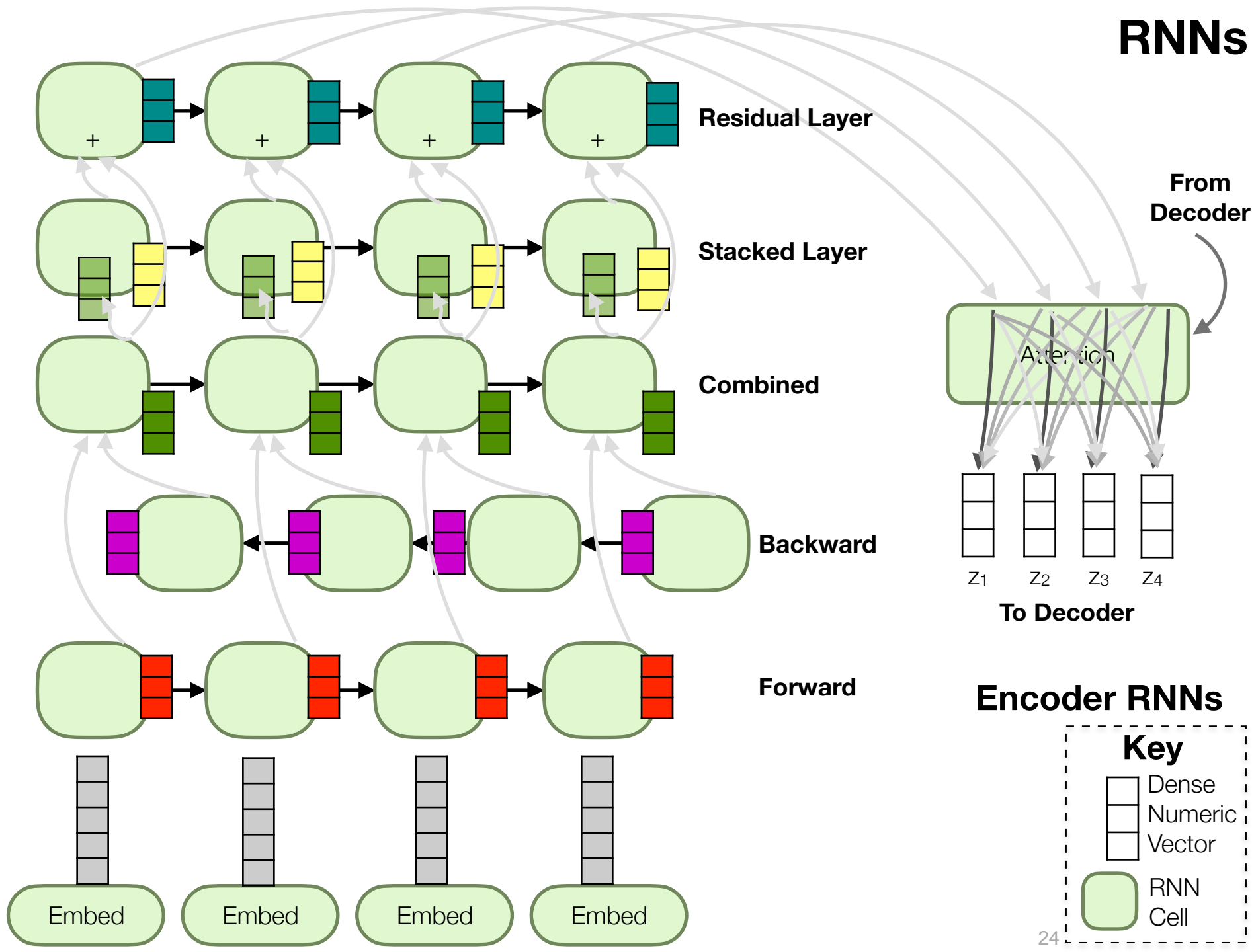
Transformers

CNN, RNN, LSTM, GAN,
Test time data,
Early stopping,
Data augmentation,
Dropout, Batch norm,
Gradient clipping

Attention

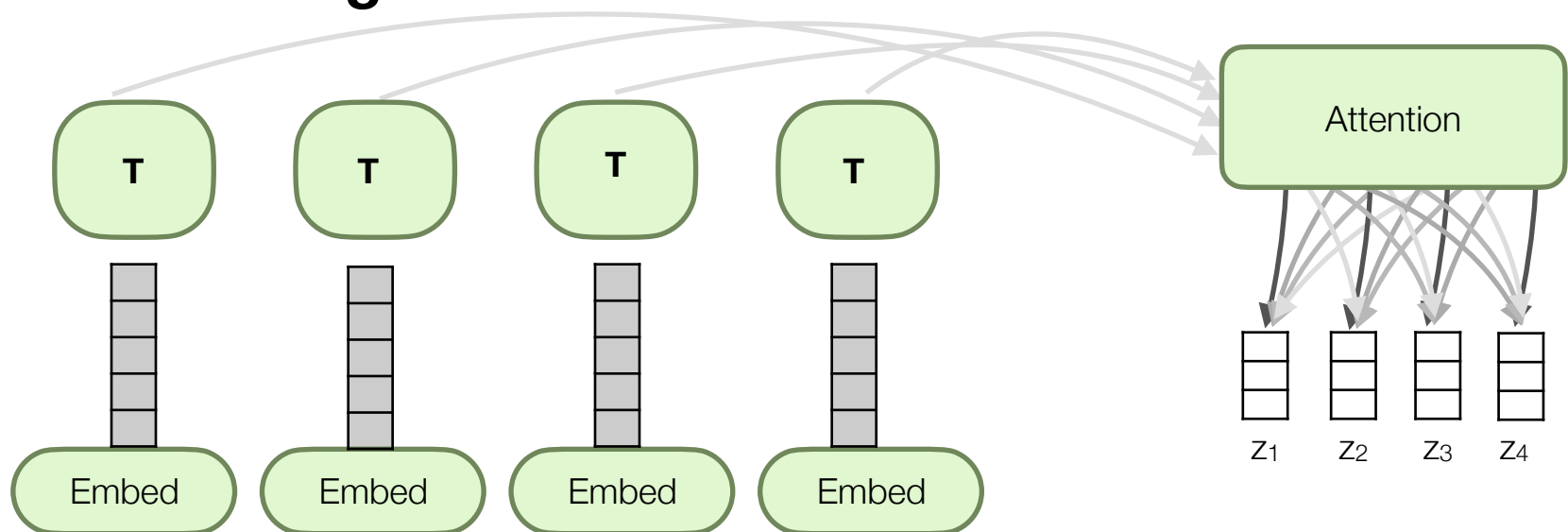


RNNs



Transformers Intuition

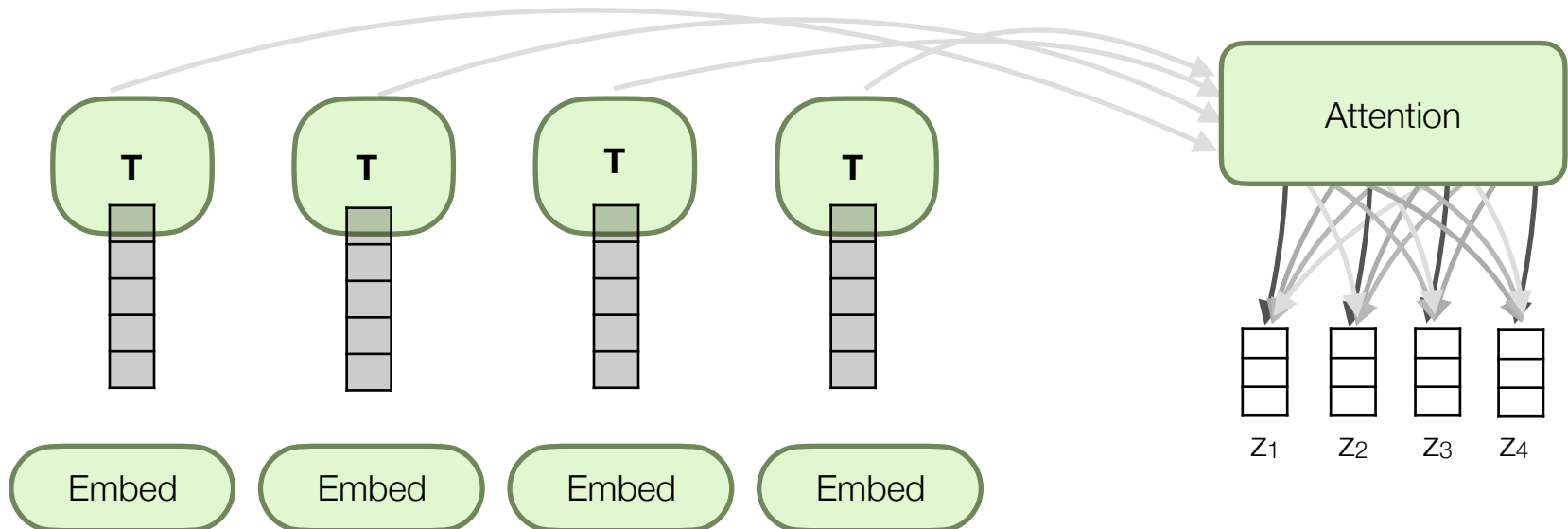
- Recurrent networks track use an “updatable” state vector, but this takes lots of iterative processing across sequence
- Attention mechanism (in RNNs) already takes a weighted sum of state vectors to generate new token in a decoder
- ... so why not just use attention on a transformation of the embedding vectors? **Do away with the recurrent state vector all together?**



Attention is All You Need

- **Transformer Solution:**

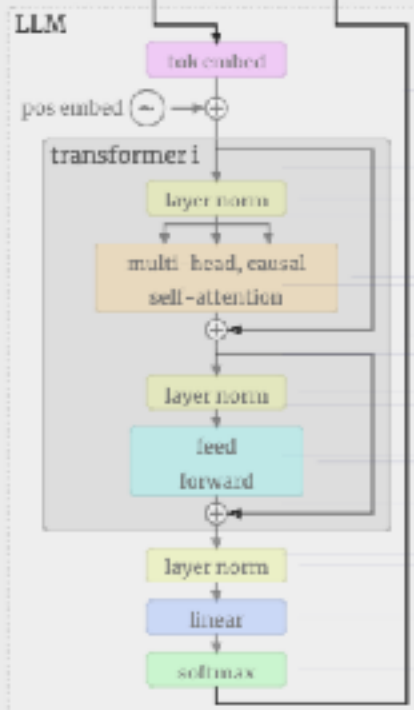
- Build attention into model from the **beginning**
- Compare all words to each other through **self-attention**
- Define a notion of “**position**” in the sequence
- *Should capture long term relationships and be highly parallelized for GPU computing!! (**ignore memory...**)*



Transformer Overview

<https://bbycroft.net/llm>

How to predict text tokens words



Components

Embedding

Layer Norm

Self

Attention

Projection

MLP

Transformer

Softmax

Output

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

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

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

GPT-2 (small)

nano-gpt

GPT-2 (XL)

GPT-3



A

A

A

A

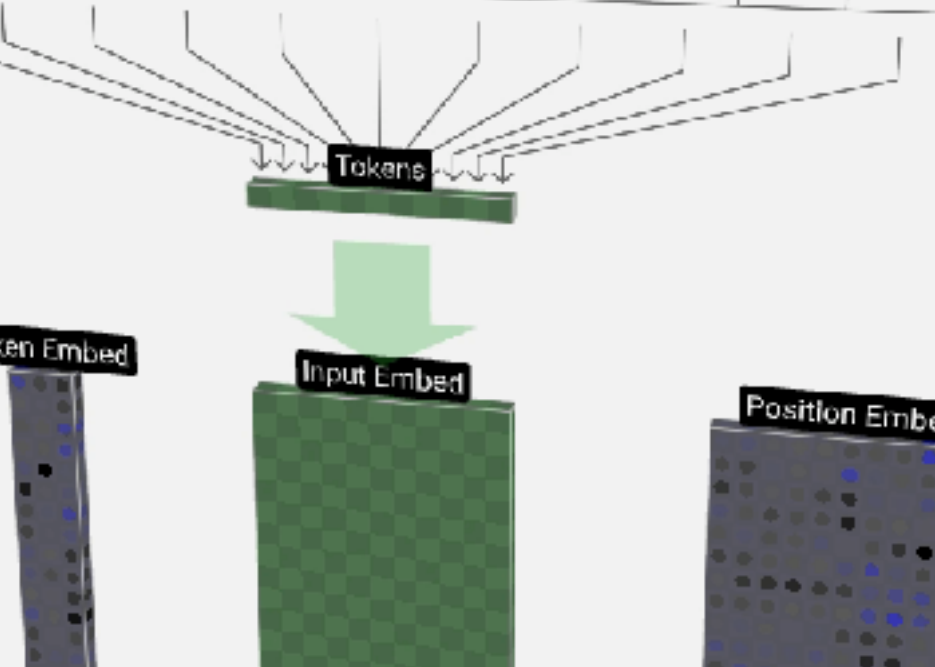
B

A

B

B

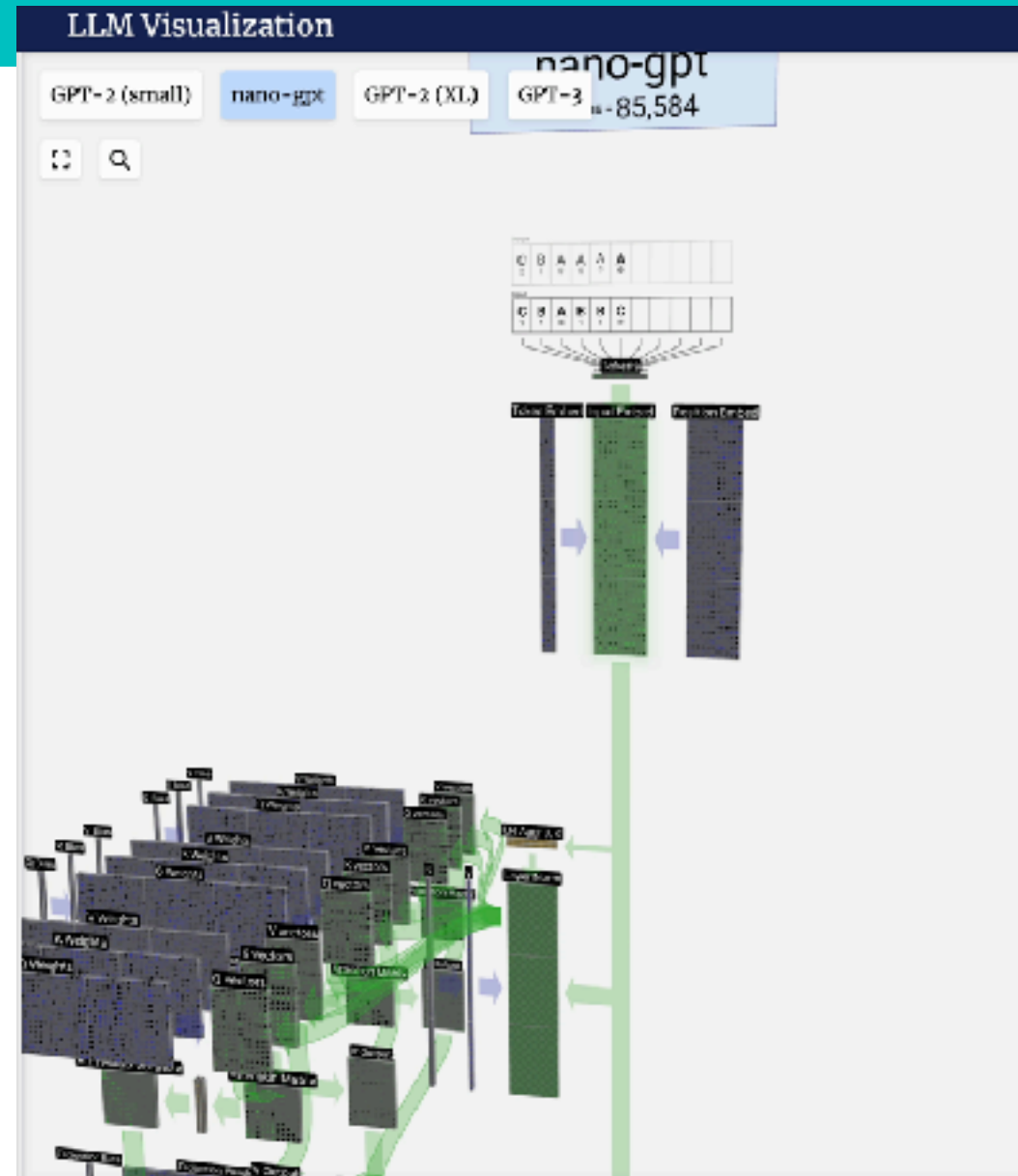
C



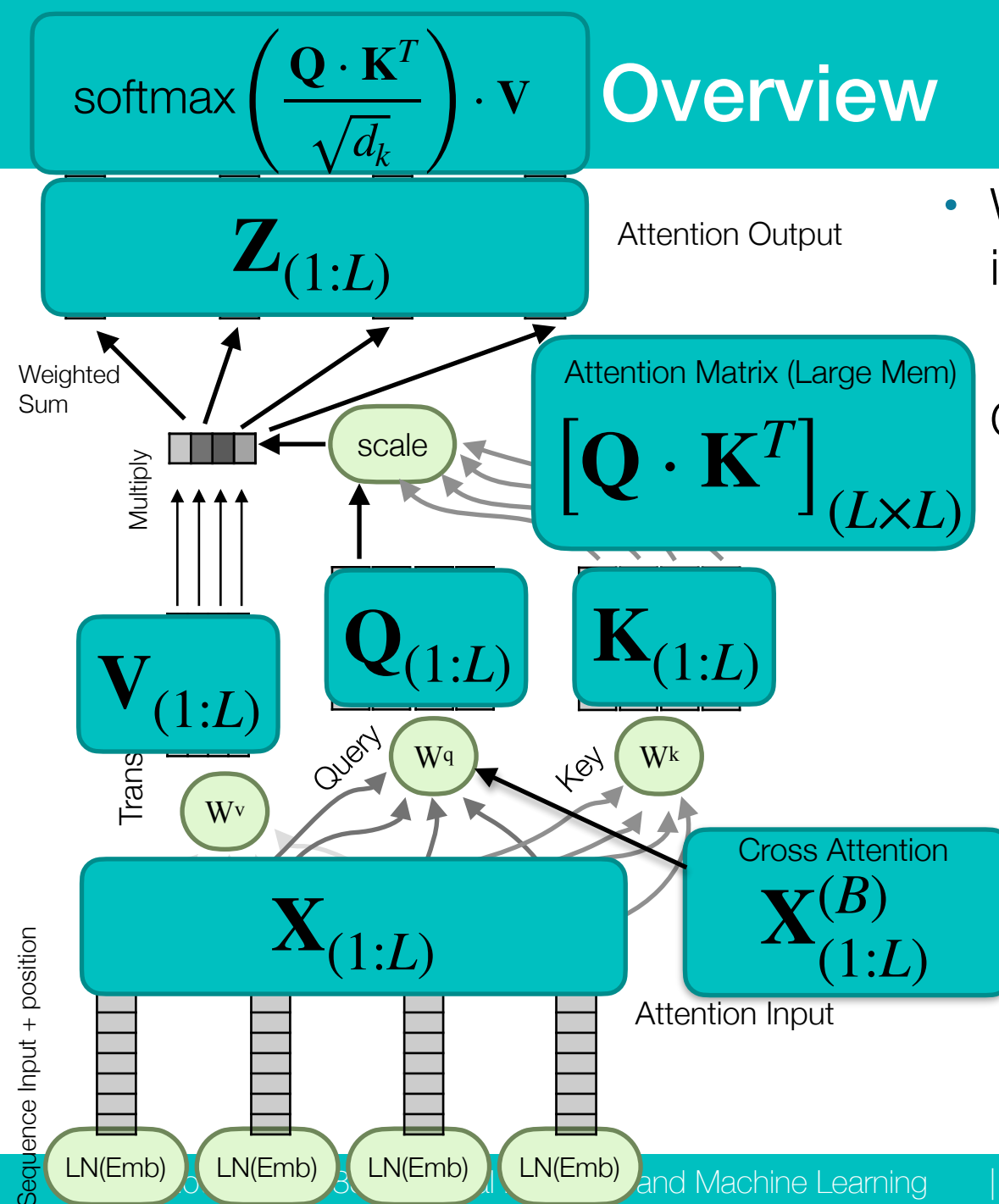
Layer Norm

Learn γ, β for each embedding column

$$LN(E^{(col)}) = \gamma_i \frac{E_i^{(col)} - \mu_E}{\sqrt{\sigma_E^2 + \epsilon}} + \beta_i$$



Overview



- What parameters are trained in diagram?

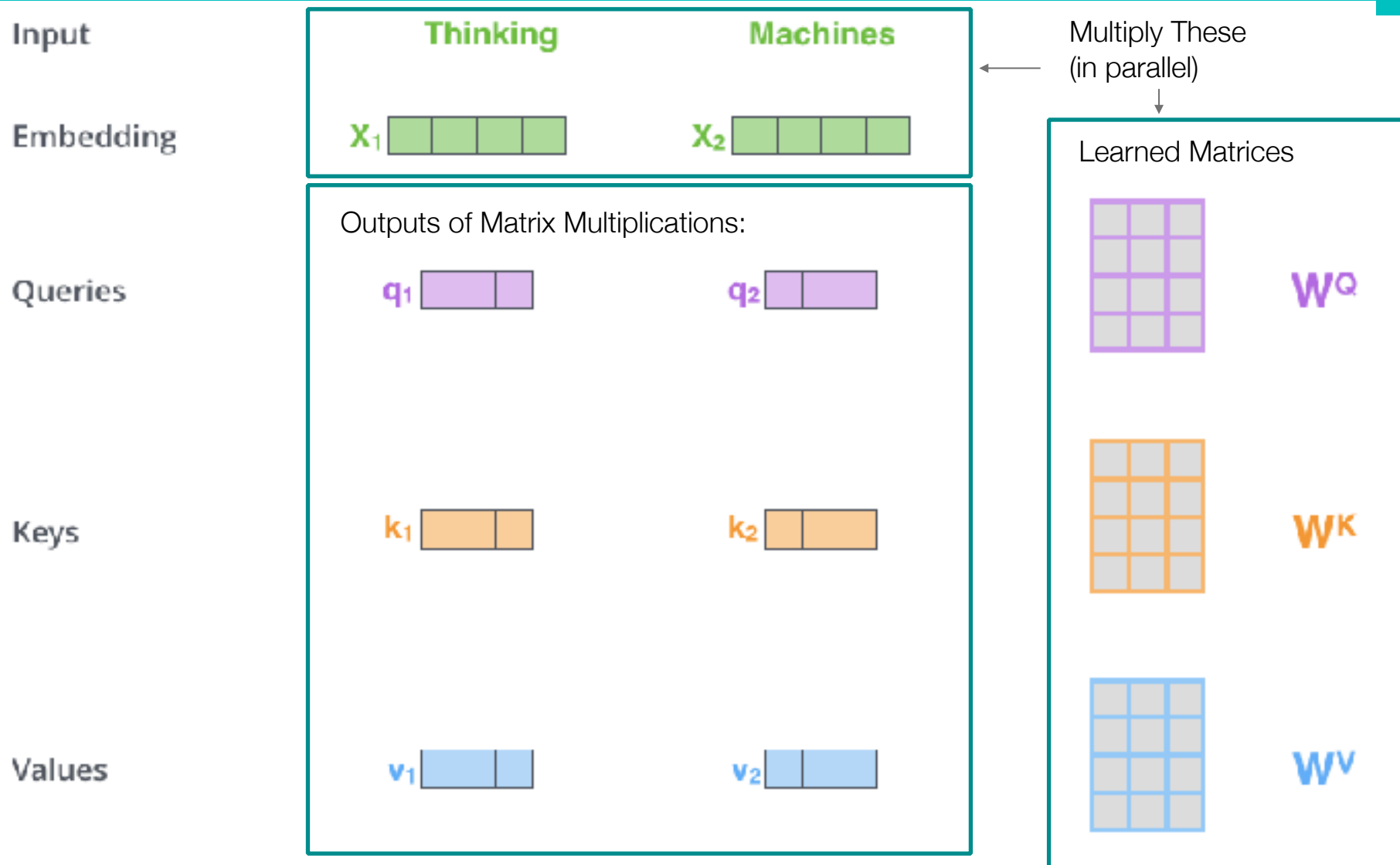
- W^v, W^q, W^k

Other Parameters:

- L : length of sequence
- Query/Key dimension, d_k
- Value dimension, d_v
- Type of positional encoding (more later)
- Cross attention versus self attention



Transformer: in more detail



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

30



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

14

Divide

12

Softmax

0.88

0.12

Multiply

Calc weights

Multiply

v_1

v_2

weighted sum for all words in document

Sum

z_1

attention for word 1

z_2

attention for word 2

Straight forward to do this operation
in matrix form:

$$\begin{matrix} \text{Thinking} \\ \text{Machines} \end{matrix} \begin{matrix} x \\ x \end{matrix} \times \begin{matrix} W^Q \\ W^Q \end{matrix} = \begin{matrix} Q \\ Q \end{matrix}$$

$\leftarrow d_k$

$$\begin{matrix} \text{Thinking} \\ \text{Machines} \end{matrix} \begin{matrix} x \\ x \end{matrix} \times \begin{matrix} W^K \\ W^K \end{matrix} = \begin{matrix} K \\ K \end{matrix}$$

$\leftarrow d_k$

$$\begin{matrix} \text{Thinking} \\ \text{Machines} \end{matrix} \begin{matrix} x \\ x \end{matrix} \times \begin{matrix} W^V \\ W^V \end{matrix} = \begin{matrix} V \\ V \end{matrix}$$

$\leftarrow d_v$

$$\text{softmax} \left(\frac{Q \times K^T}{\sqrt{d_k}} \right) \begin{matrix} V \\ V \end{matrix}$$

$$= \begin{matrix} Z \\ Z \end{matrix}$$

z_1
 z_2

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

Size of Q,K,V:
 $|\text{Seq Len}| \times d_v$ or d_k

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

Self Attention: From <https://>

ight forward to do this operation
atrix form:

$$\begin{matrix} \text{Thinking} \\ \text{Machines} \end{matrix} \begin{matrix} \mathbf{X} \\ \mathbf{W}^Q \end{matrix} = \mathbf{Q}$$

$$\begin{matrix} \text{Thinking} \\ \text{Machines} \end{matrix} \begin{matrix} \mathbf{X} \\ \mathbf{W}^K \end{matrix} = \mathbf{K}$$

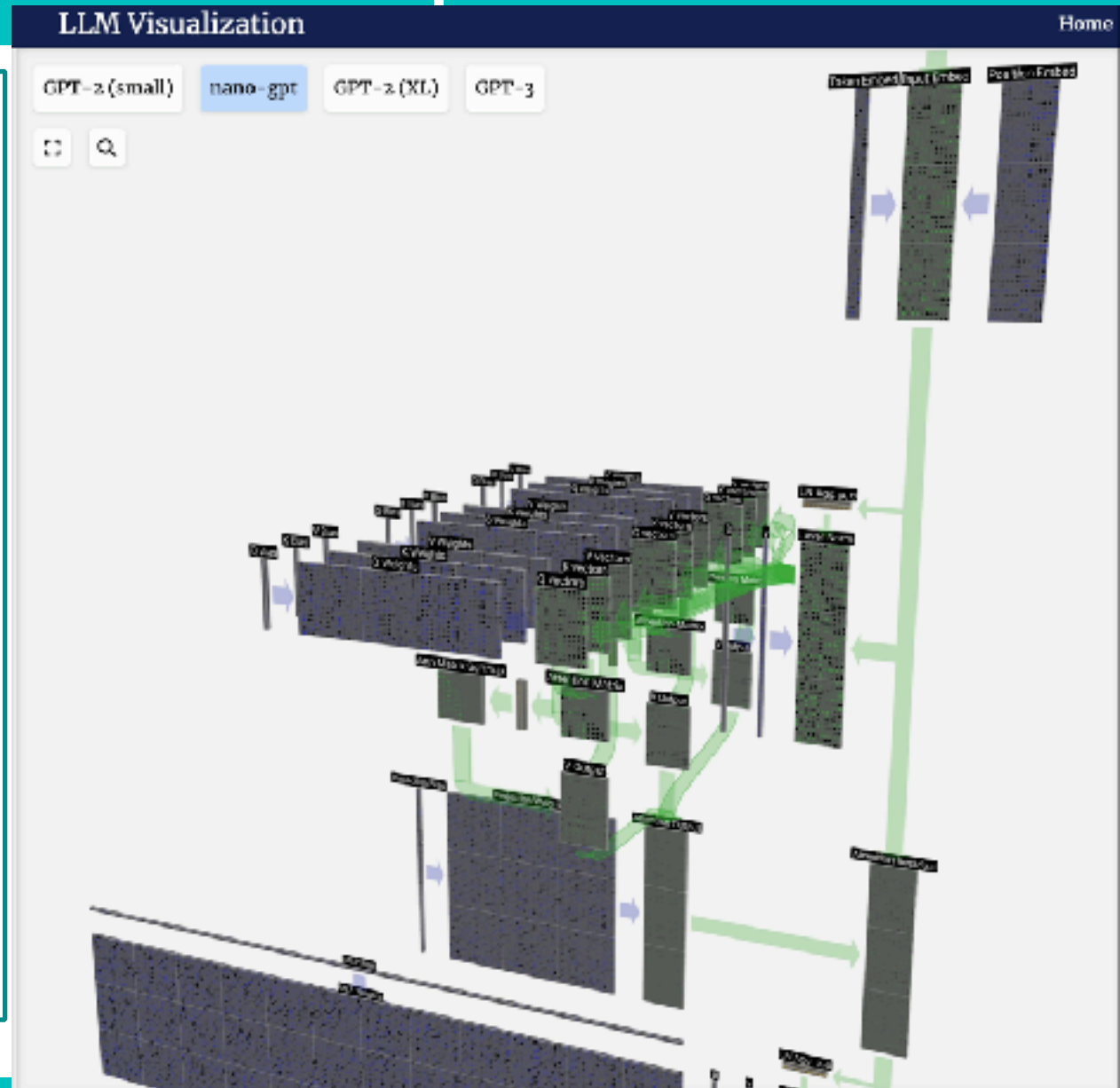
$$\begin{matrix} \text{Thinking} \\ \text{Machines} \end{matrix} \begin{matrix} \mathbf{X} \\ \mathbf{W}^V \end{matrix} = \mathbf{V}$$

$$\text{ftmax} \left(\frac{\mathbf{Q} \times \mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V}$$

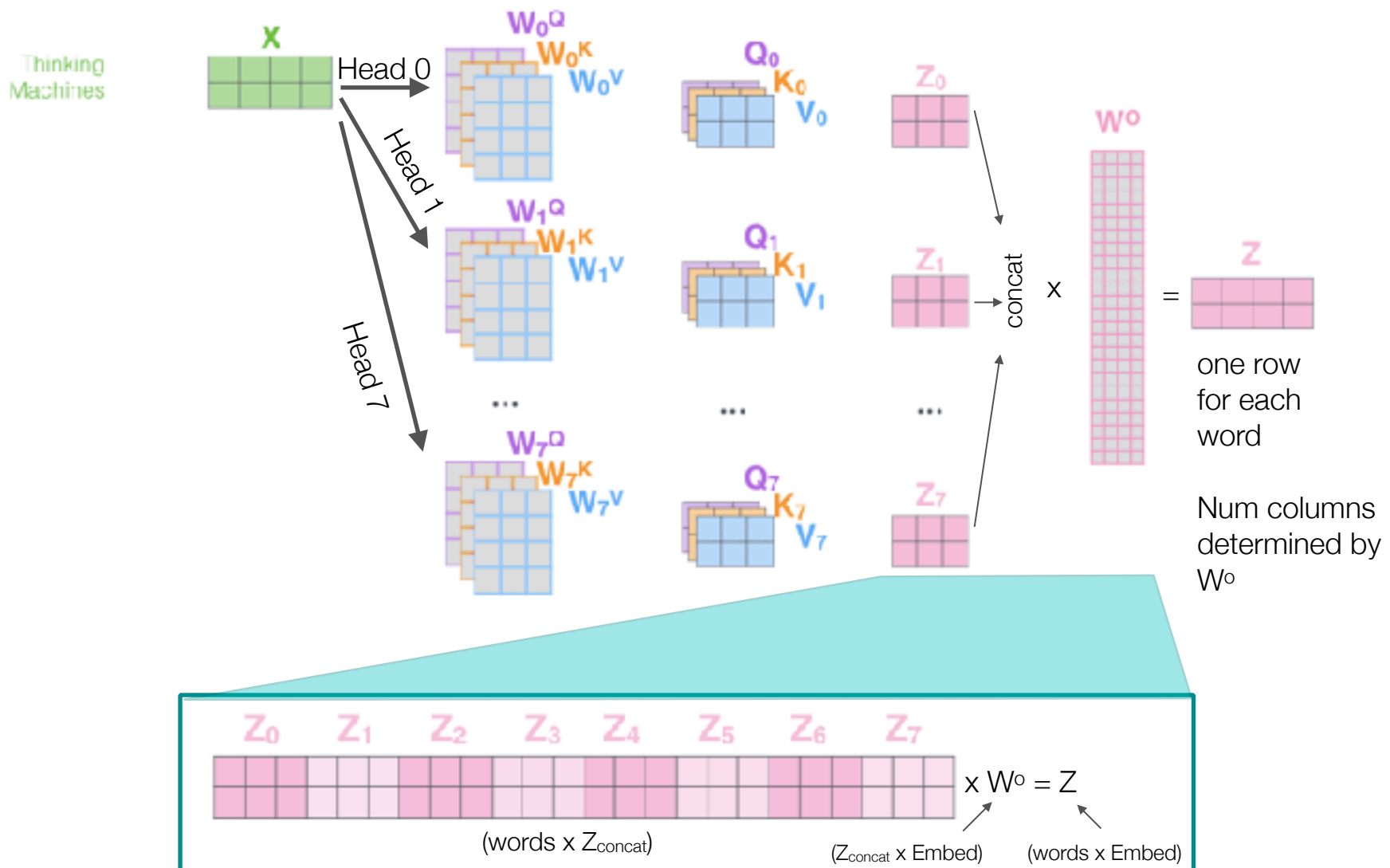
$$\mathbf{Z} = \begin{matrix} \mathbf{Z}_1 \\ \mathbf{Z}_2 \end{matrix}$$

Can perform more than once!
Multiple heads!

output of each head is \mathbf{Z}_i



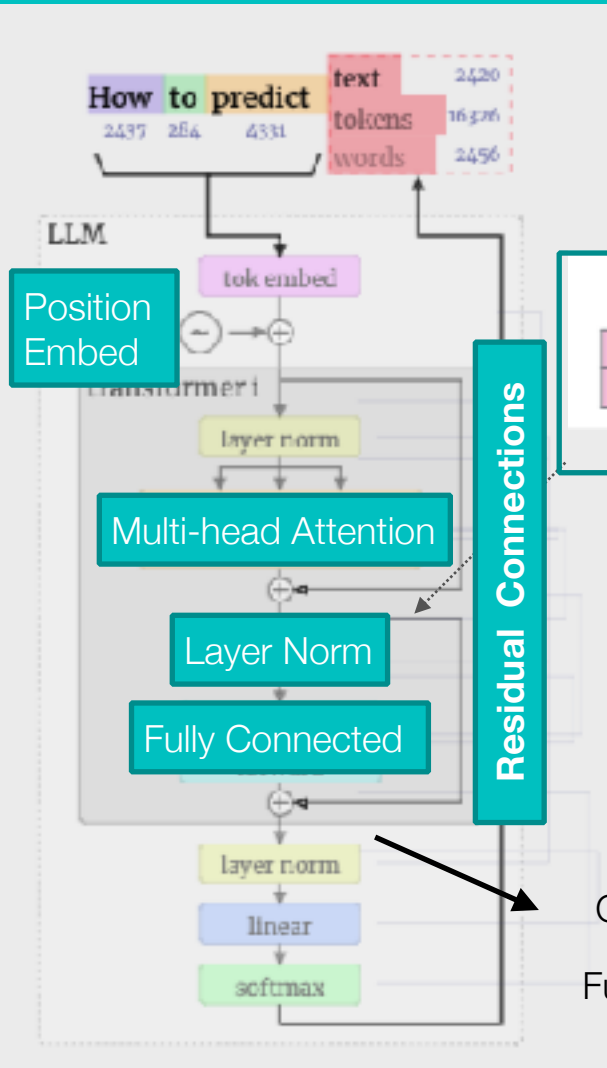
Transformer: Multi-headed Attention



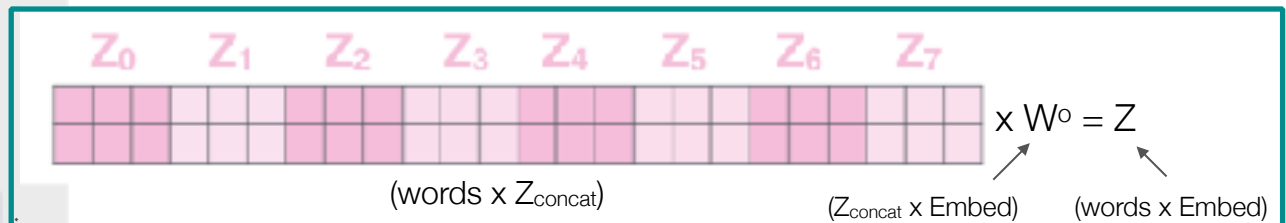
Excellent Blog on Transformers: <http://jalammar.github.io/illustrated-transformer/>



Putting It Together



```
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, dims of  $W^0$ )
```



$$LN(Z_{row}) = \gamma \frac{Z_i - \mu_Z}{\sqrt{\sigma_Z^2 + \epsilon}} + \beta$$

