

Lecture Notes for **Machine Learning in Python**



Professor Eric Larson
Sequential CNN and Transformers

Lecture Agenda

- Logistics
 - Grading Update
 - Lectures
 - Sequential Networks due **during finals**
- Agenda
 - CNNs for Sequential Processing (review)
 - Transformers

Class Overview, by topic

Table Data
Visualization

Numpy, Pandas, Seaborn
Overviews with some in-depth discussion

Dimension
Reduction and
Image Processing

Scikit-learn, Scikit Image,
Intuition only, Some mathematics

Linear and
Logistic
Regression

Numpy, Recreate API for Scikit-learn
Detailed mathematics for simple optimization
intuition for advanced optimization

Neural Networks
and Back Prop.

Numpy
Detailed mathematics for NN operations

Wide and Deep
Networks

Convolutional
Networks

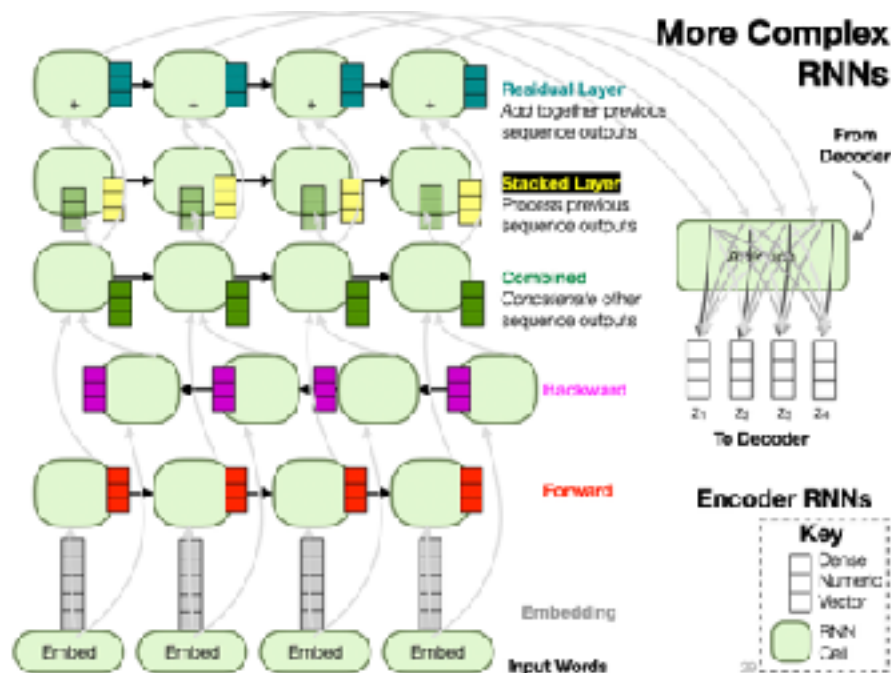
Sequential
Networks

Keras, Tensorflow
Intuition, Detailed implement.

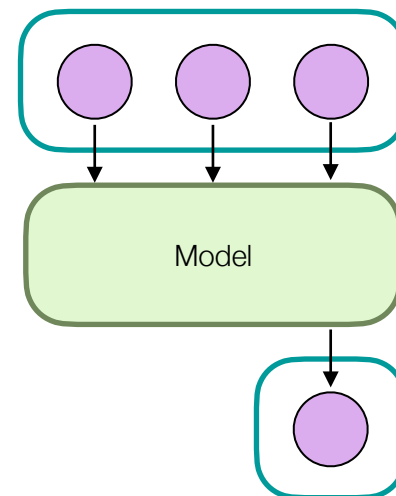
Ethics in
Language Models

ConceptNet
Case studies

Last Time:



Many to One

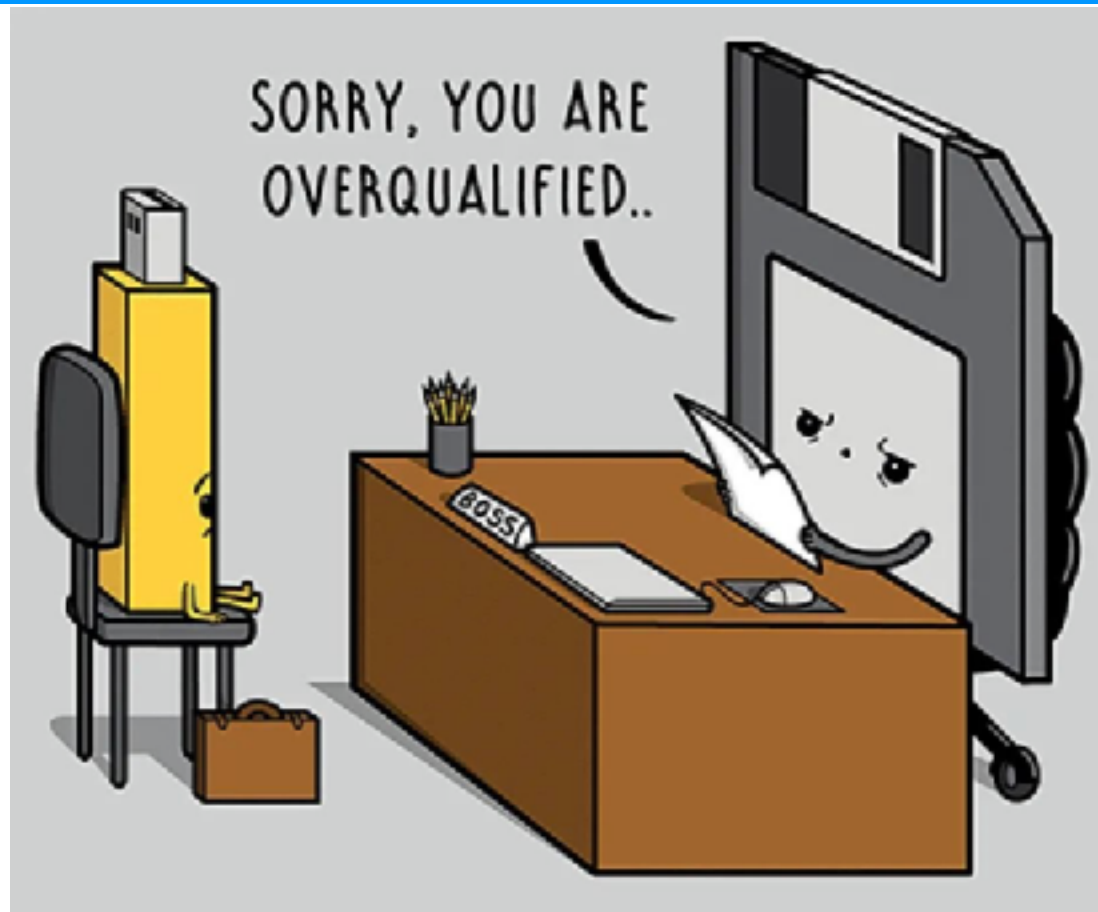


Visualization

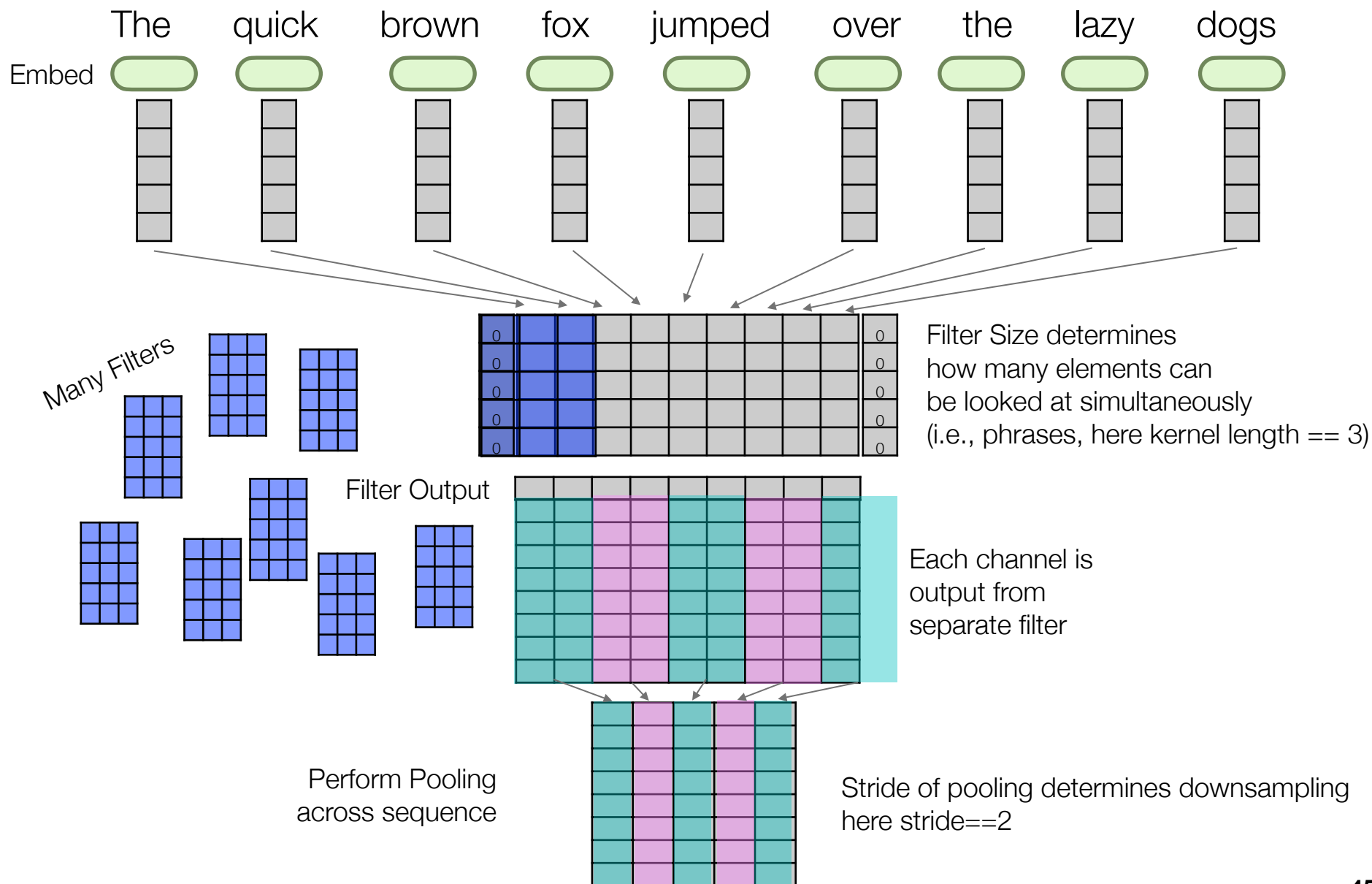
GloVe produces word vectors with a marked banded structure that is evident upon visualization:



CNNs for Sequences

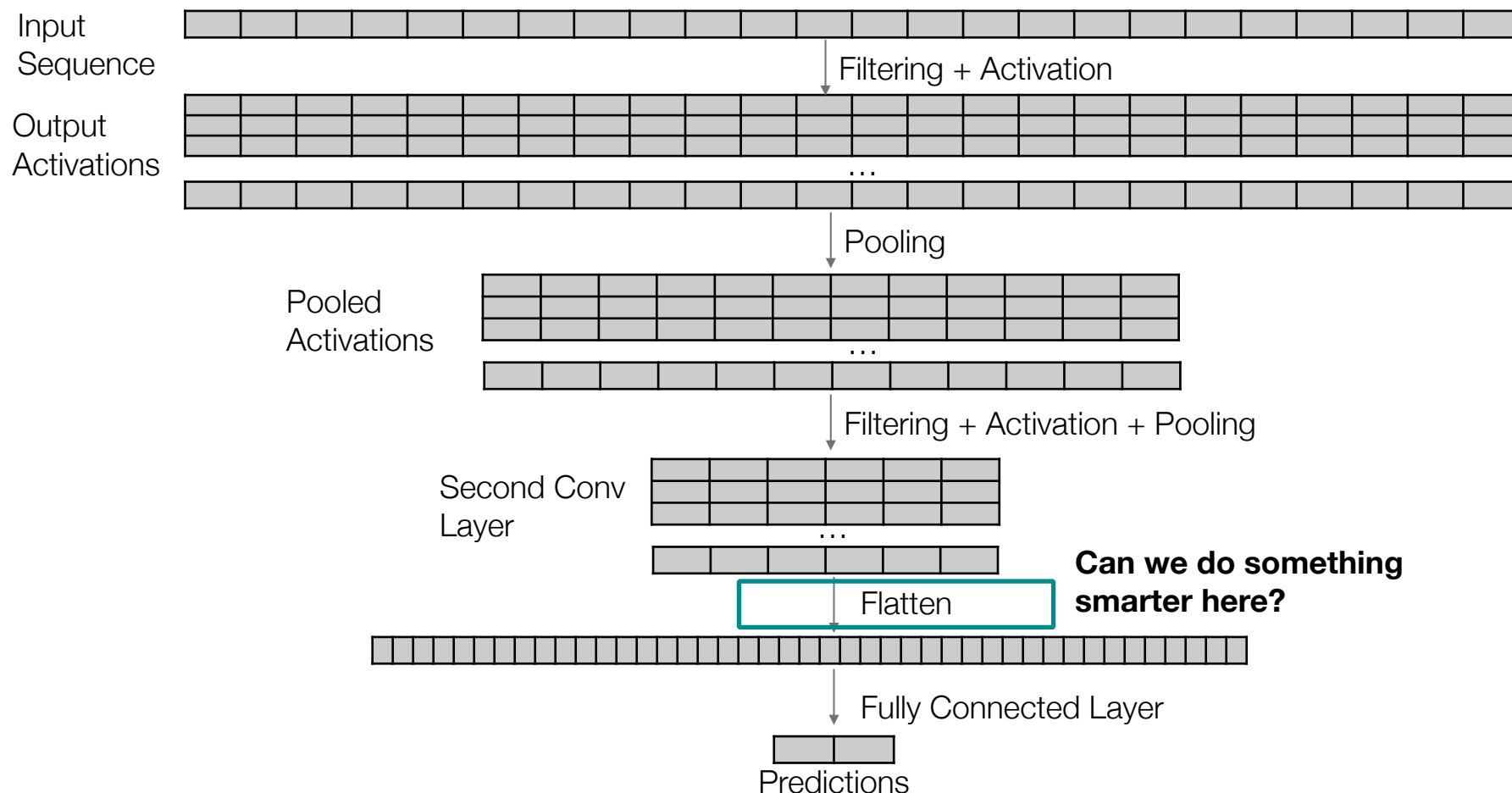


CNNs for Sequences



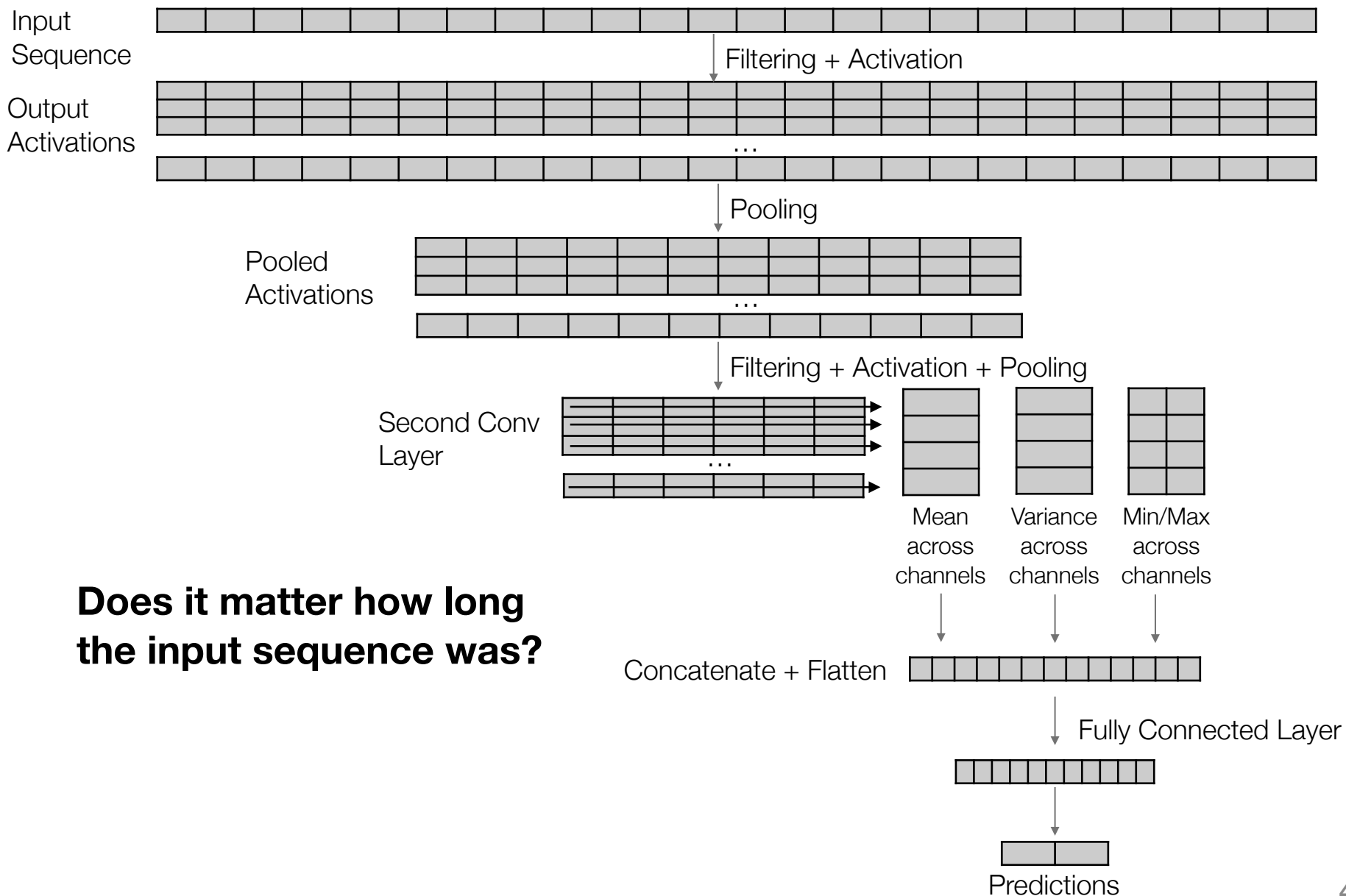
CNNs for Sequences

- RNNs are not inherently parallelized or efficient at remembering based on state vector, but CNNs can be run in parallel groups

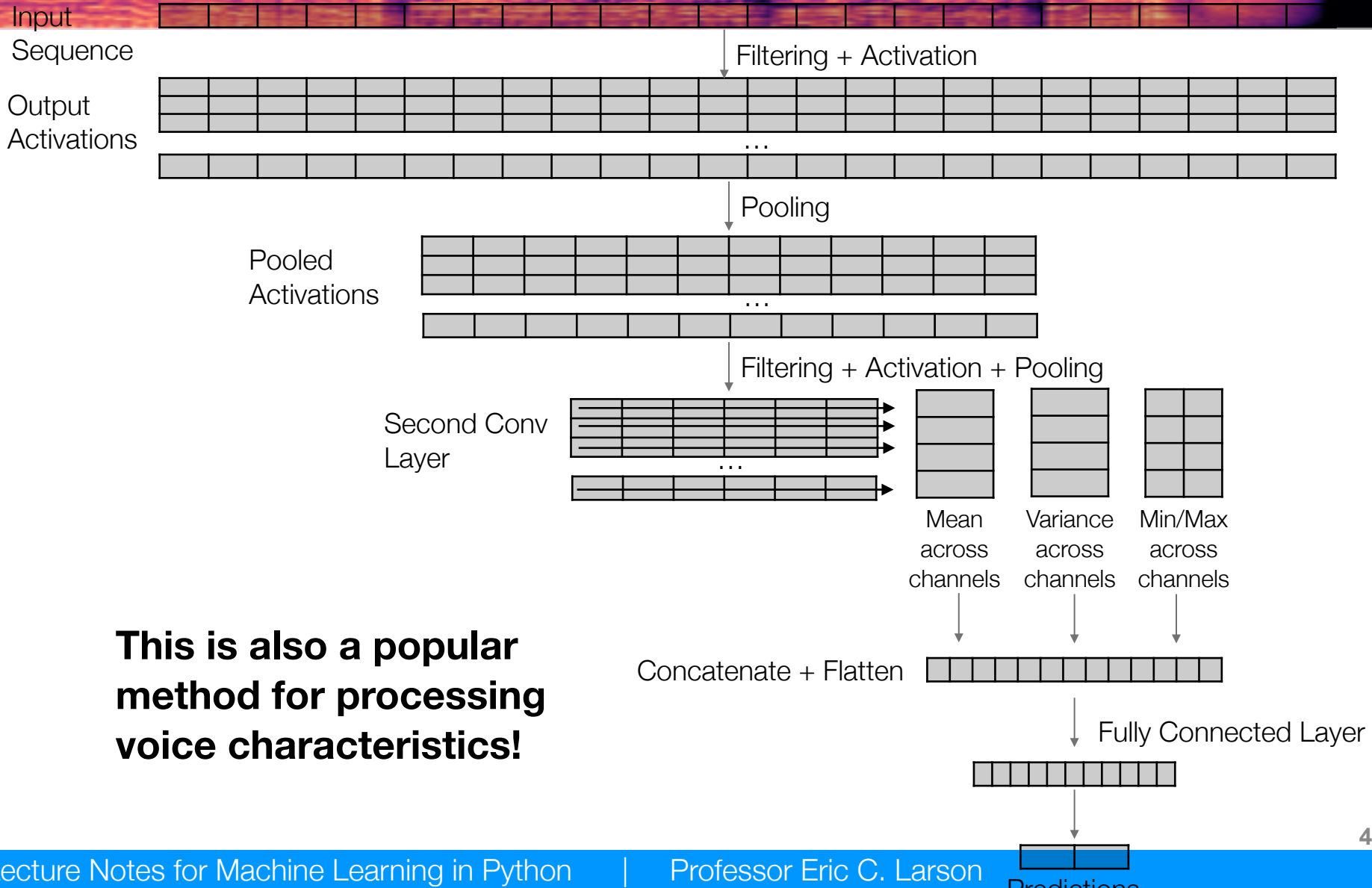


- Everything we learned in 2D CNNs can be applied to 1D CNNs...
- Residuals, separable convolution, squeezing, everything

CNNs for Sequences



CNNs for Sequences



The Sequential CNN
IMdB sentiment analysis

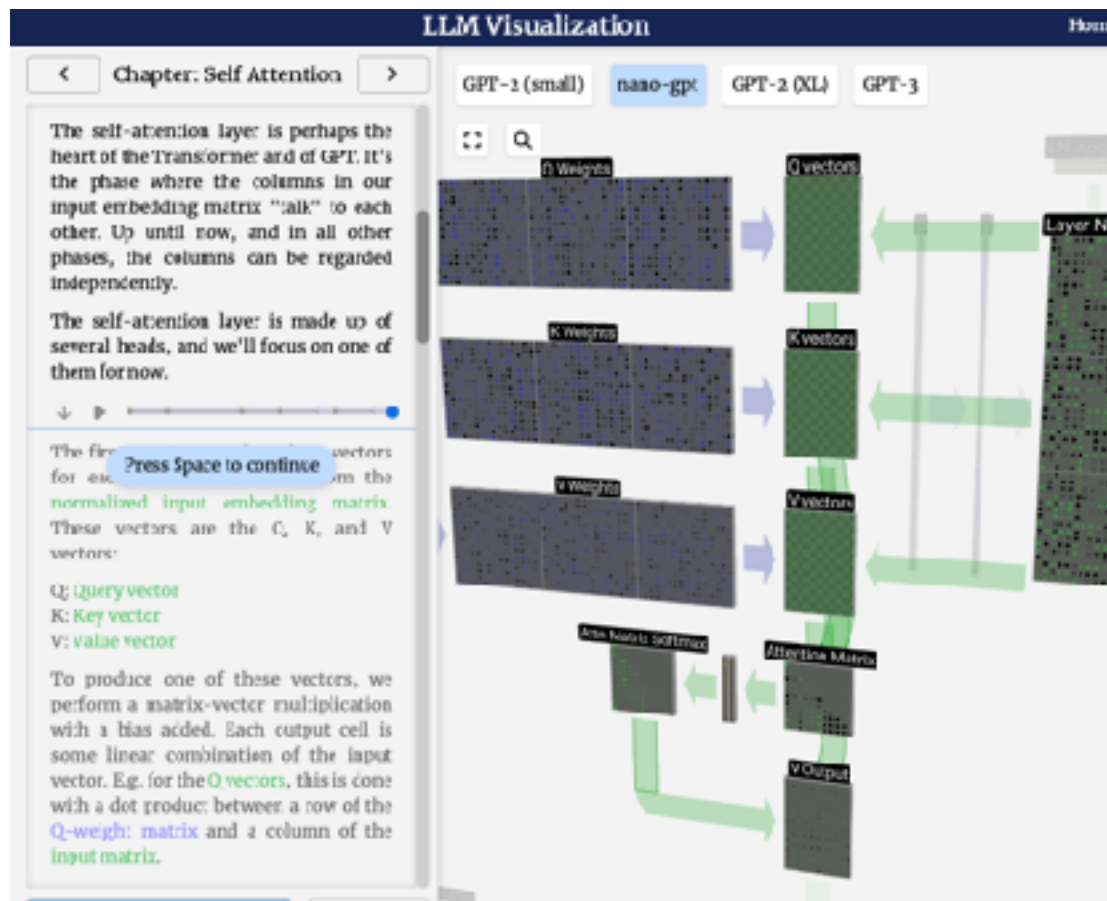


13a. Sequence Basics [Experimental].ipynb

Marked as “Experimental” because it does not yet have a textbook definition of transformers.

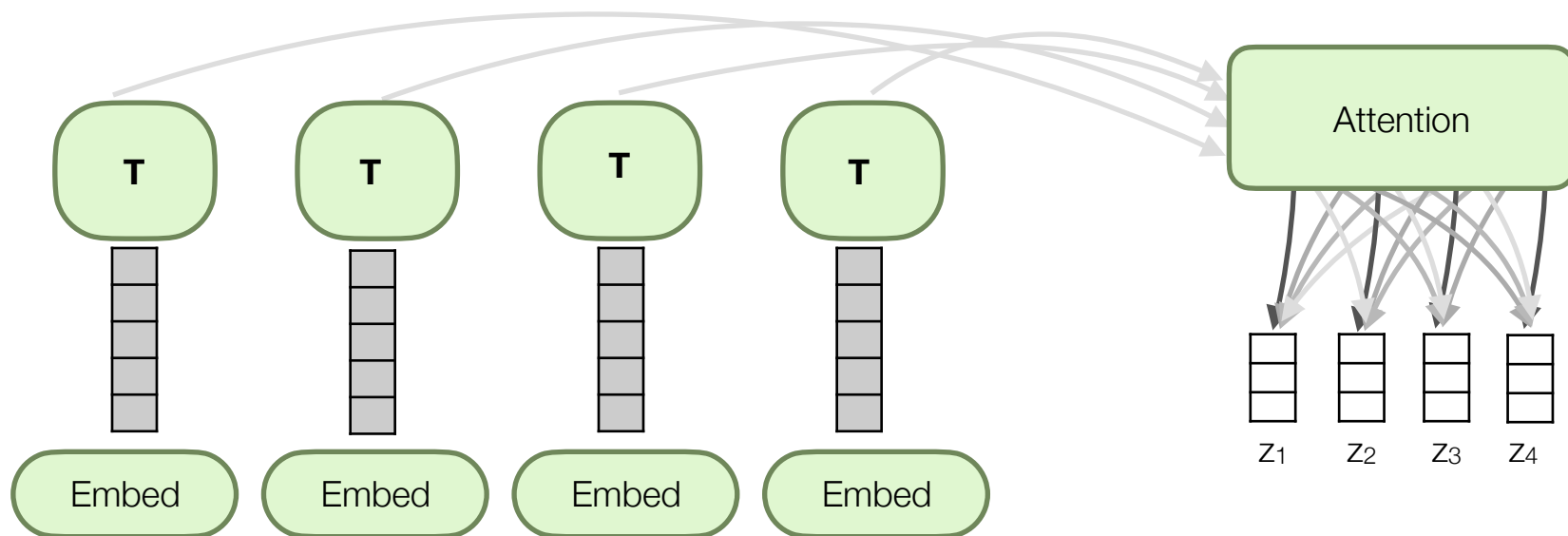
This link is perhaps the greatest tutorial on
X-formers I have ever seen

<https://bbycroft.net/llm> **Transformers**



Transformers Intuition (reminder)

- Recurrent networks track state using an “updatable” state vector, but this takes processing iterative
- 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?**



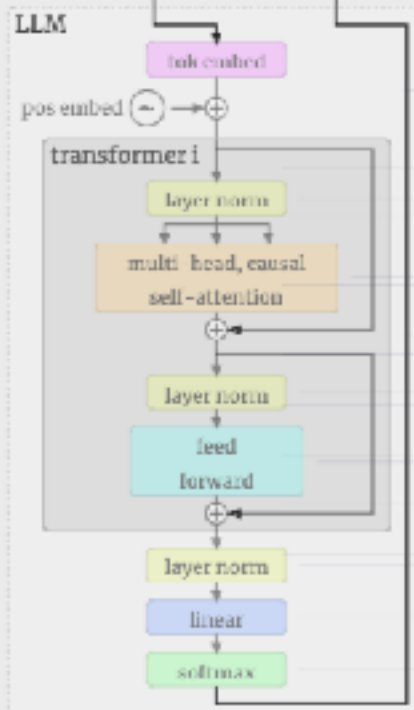
Attention is All You Need

- **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-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

<https://bbycroft.net/llm>

How to predict text
 14,37 284 4,133
 tokens 14,376
 words 24,568



Components

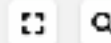
Embedding
 Layer Norm
 Self Attention
 Projection
 MLP
 Transformer
 Softmax
 Output

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

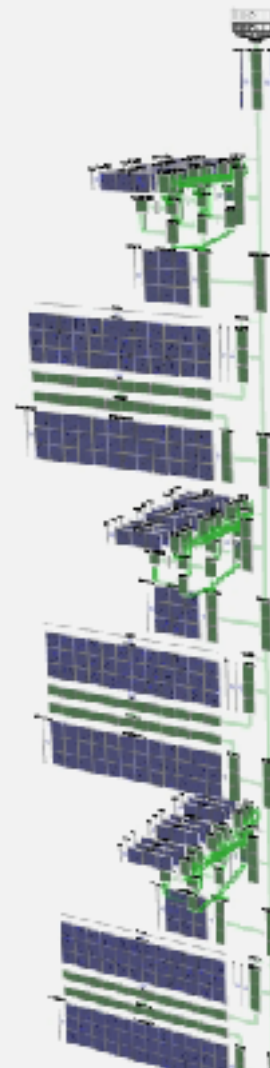
$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

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



tokens = 85,584



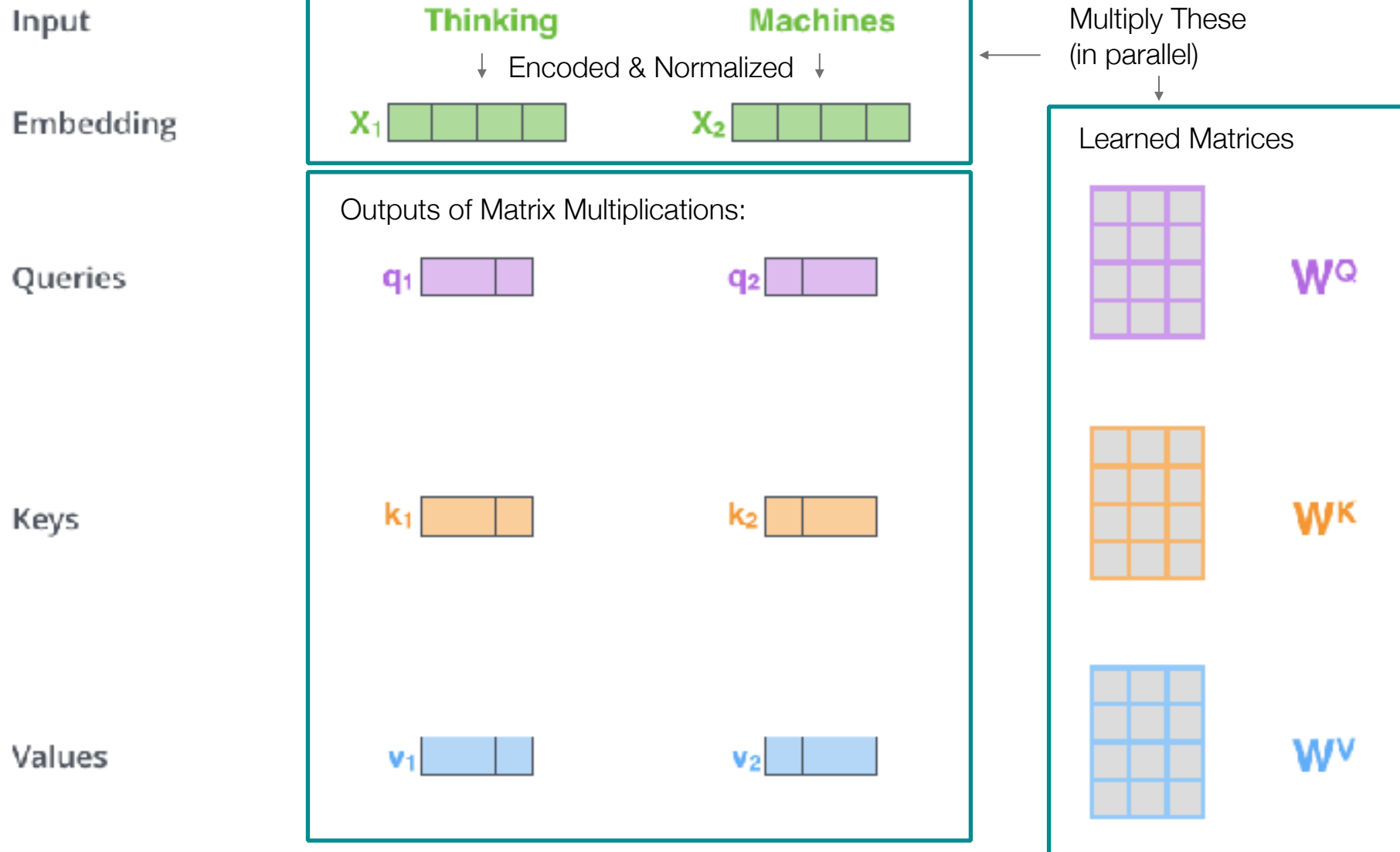
Layer Norm

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

Learn to normalize
along a dimension of Z

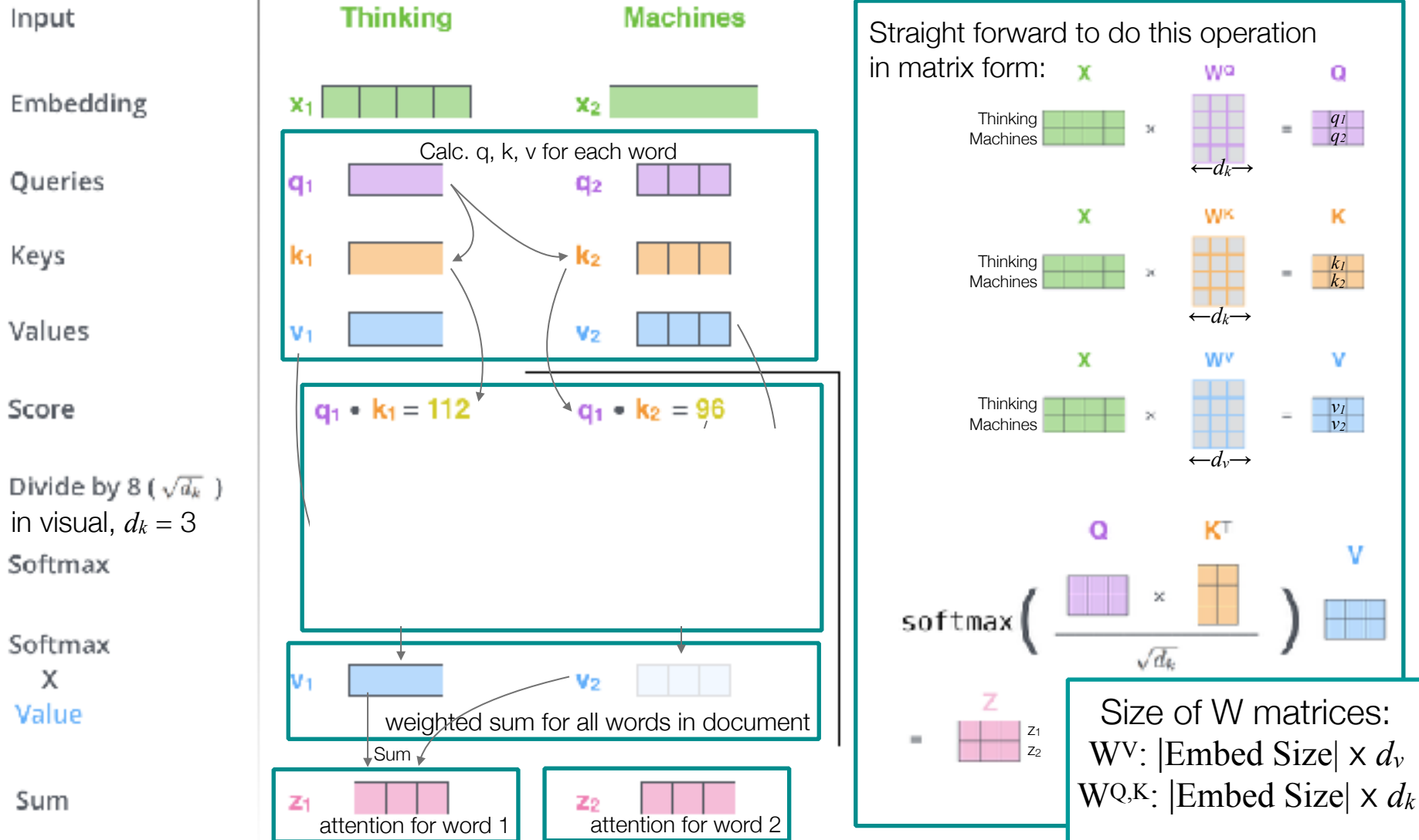


Transformer: in more detail



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

Transformer: in more detail



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

Self Attention: From <https://bbbycroft.net/llm>

right forward to do this operation
matrix form:

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

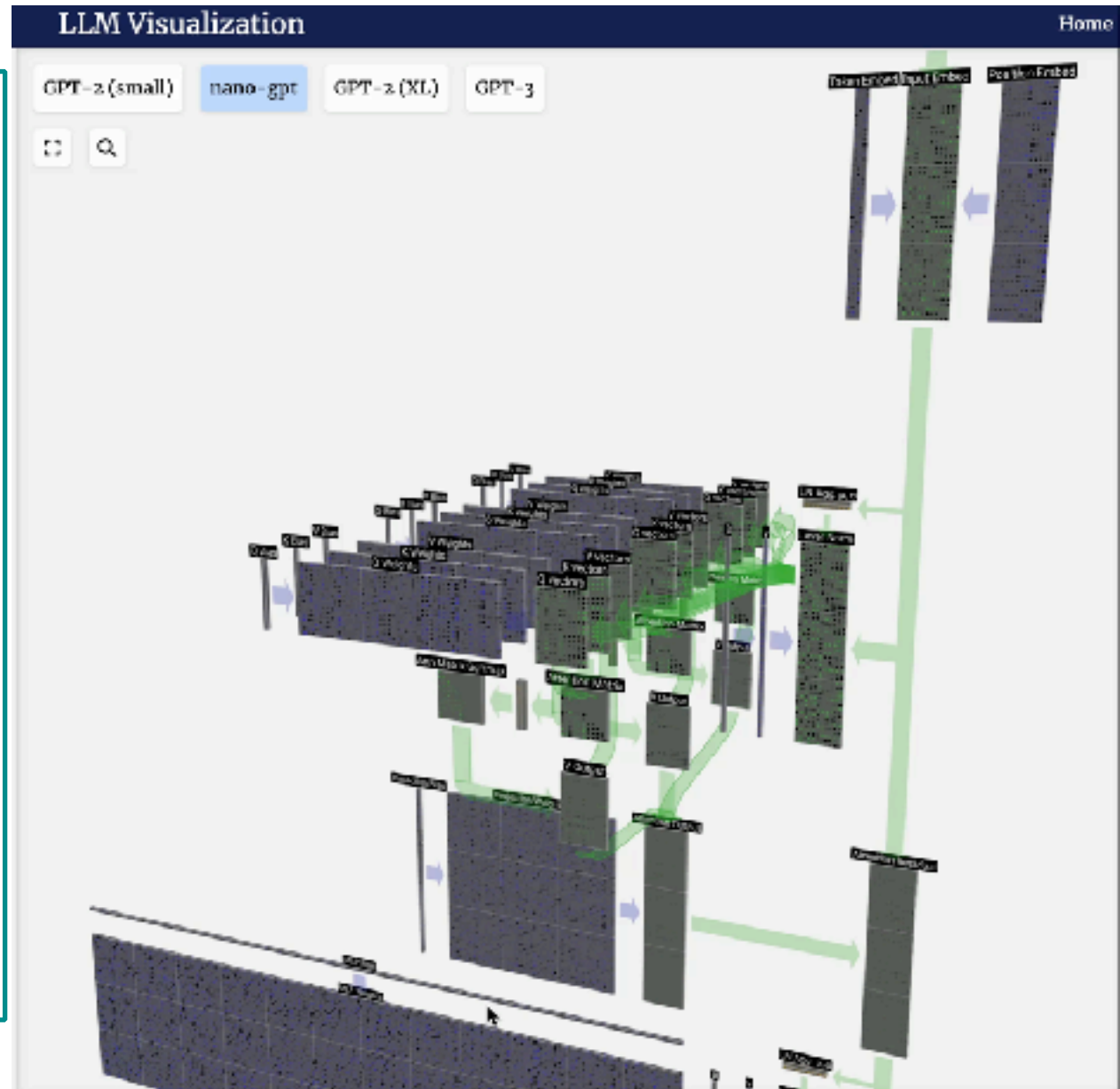
$$\begin{matrix} \text{Thinking} \\ \text{Machines} \end{matrix} \begin{matrix} \mathbf{X} \\ \mathbf{W}^K \end{matrix} = \mathbf{K}$$

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

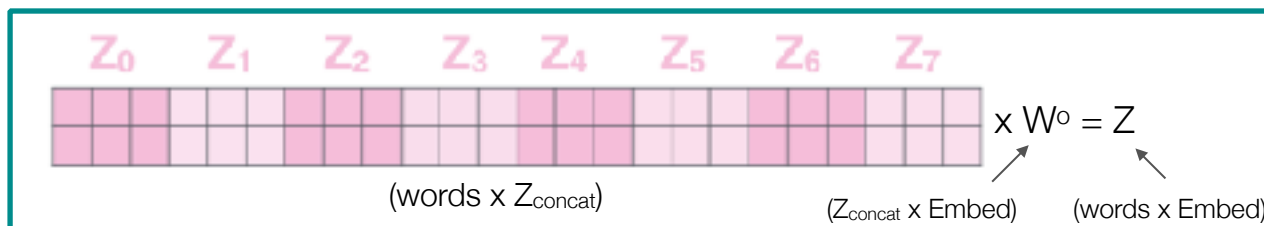
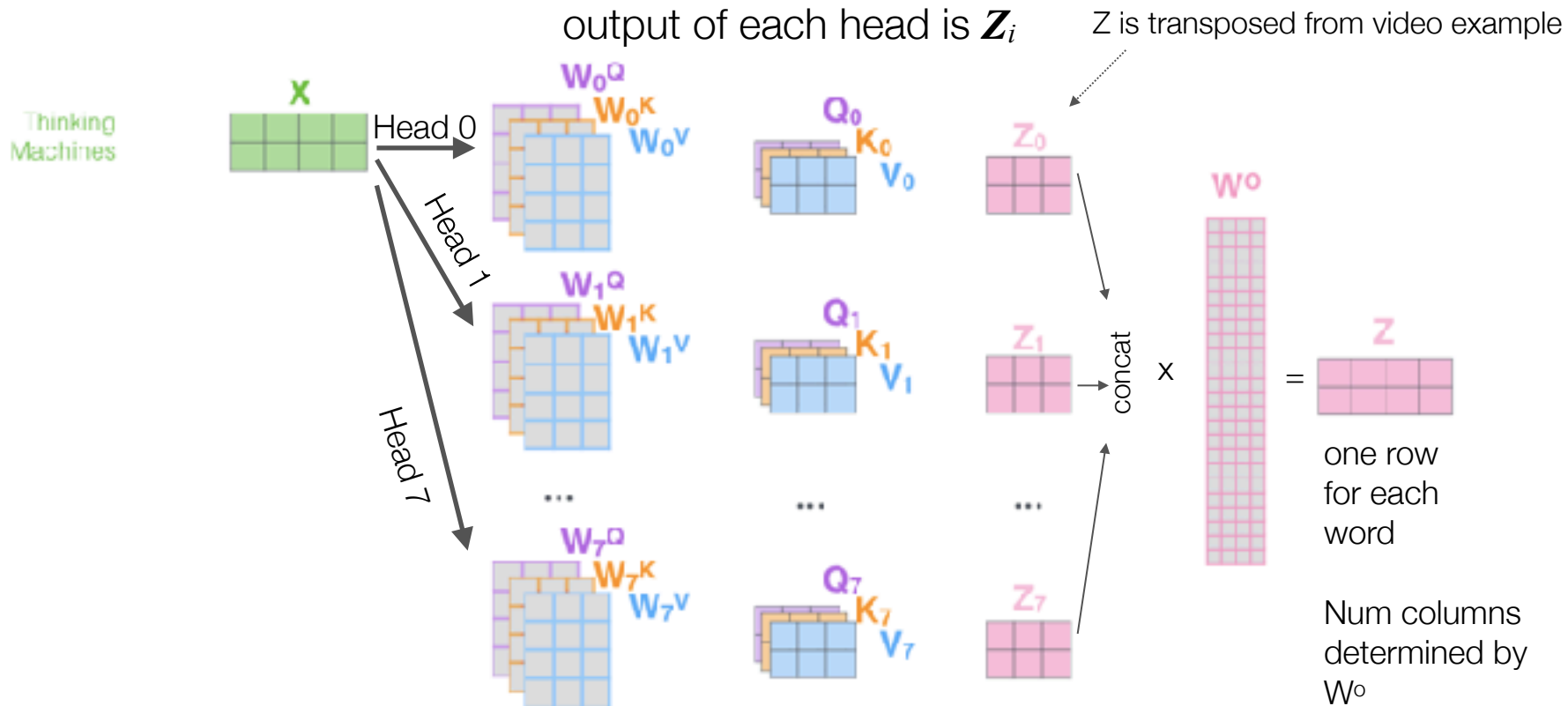
$$\text{softmax} \left(\frac{\mathbf{Q} \times \mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V}$$

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

output of each head is \mathbf{Z}_i



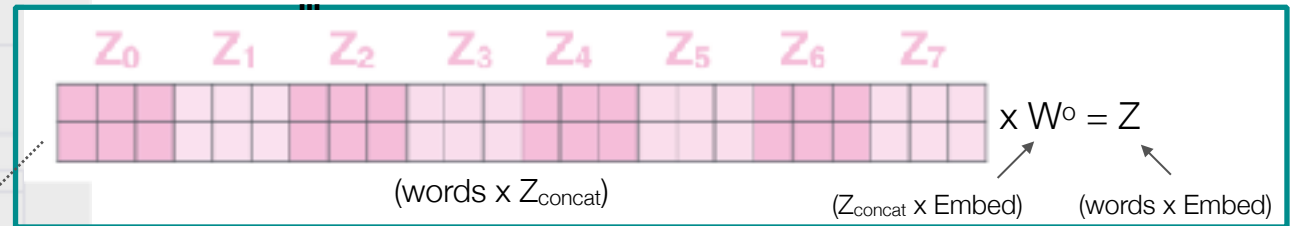
Transformer: Multi-headed Attention



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

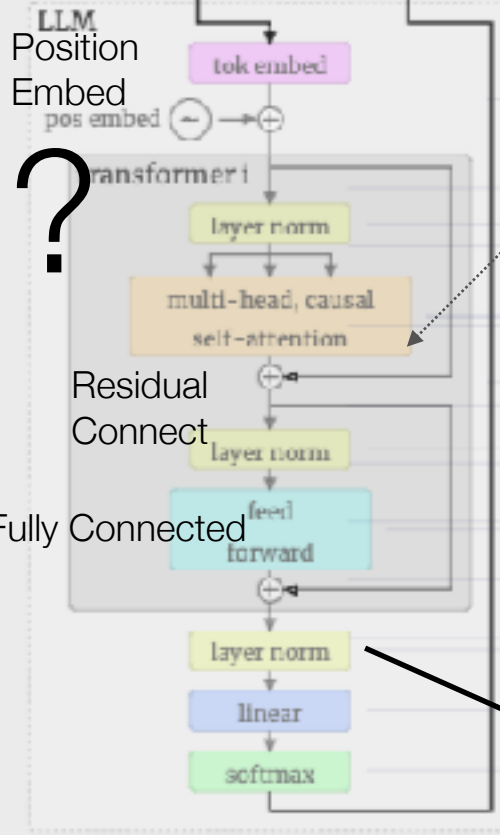
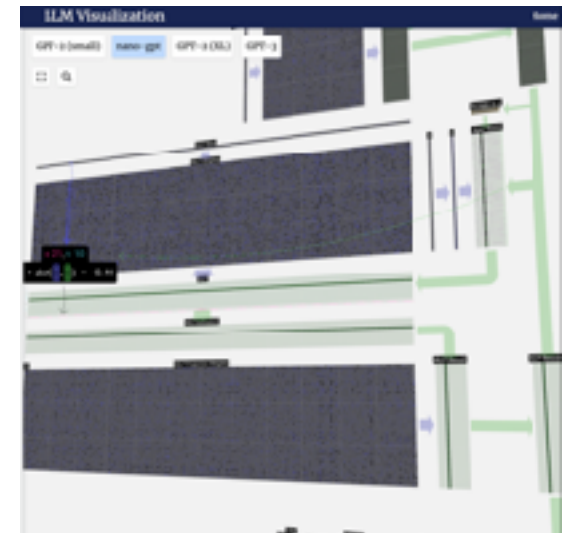
Putting It Together

```
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  $W^o$ )
```



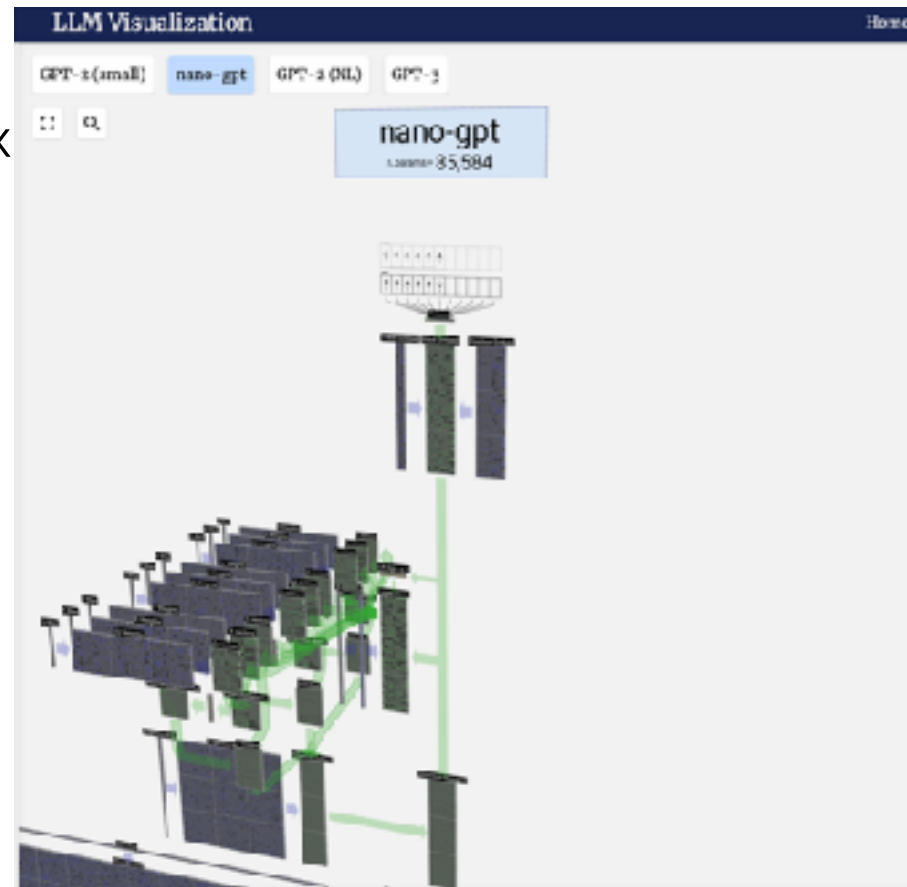
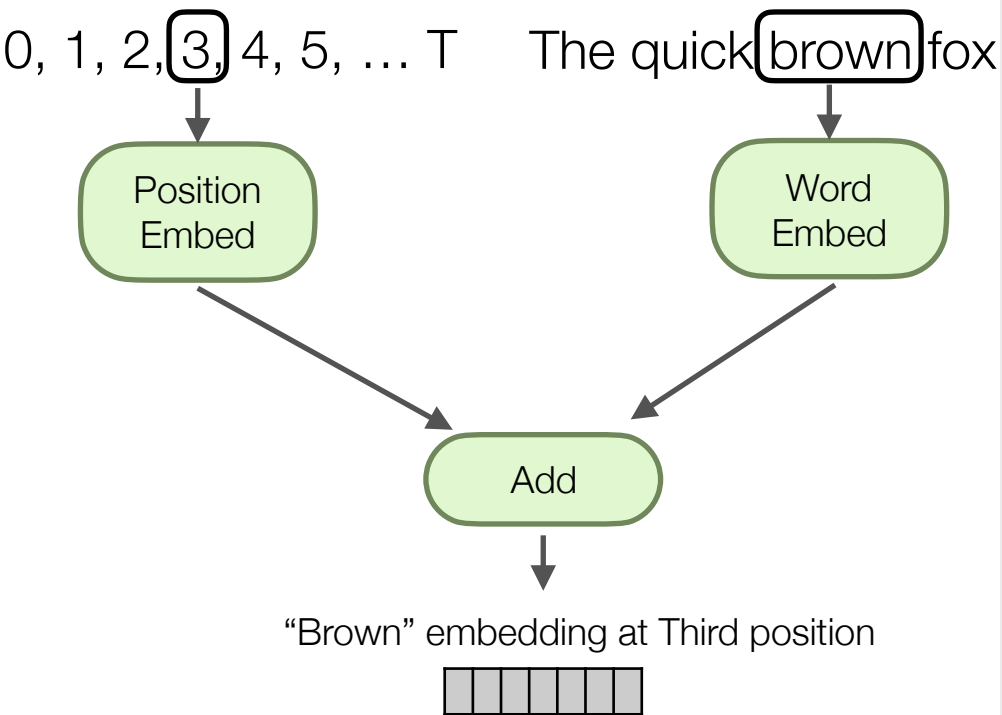
$$LN(Z_{row}) = \gamma \frac{Z_i - \mu_Z}{\sqrt{\sigma_Z^2 + \epsilon}} + \beta$$

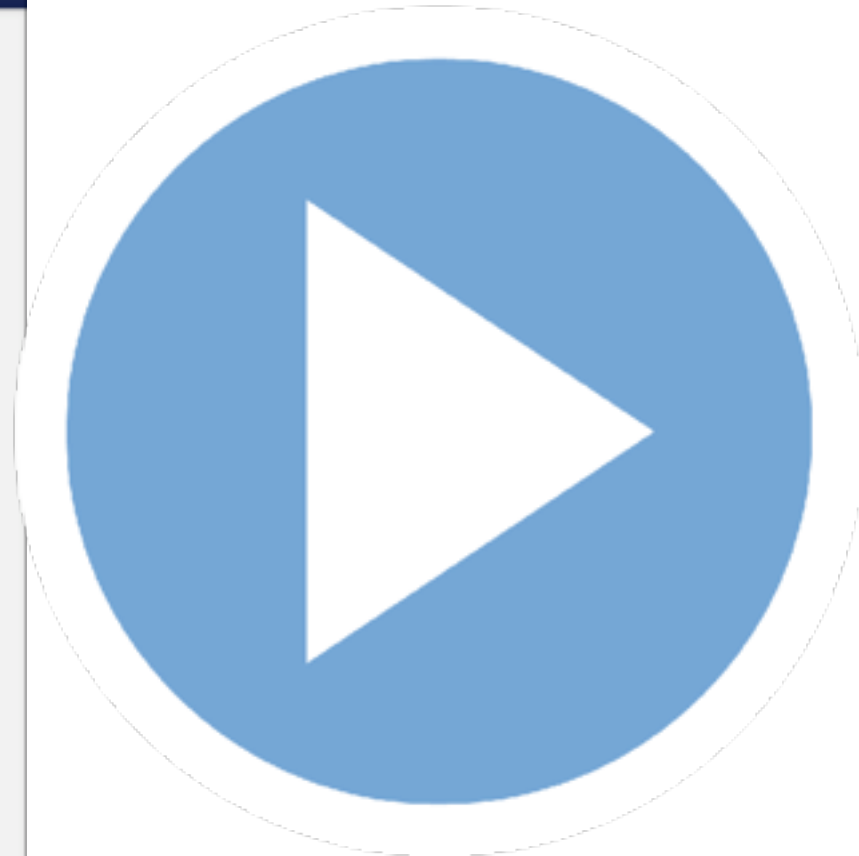
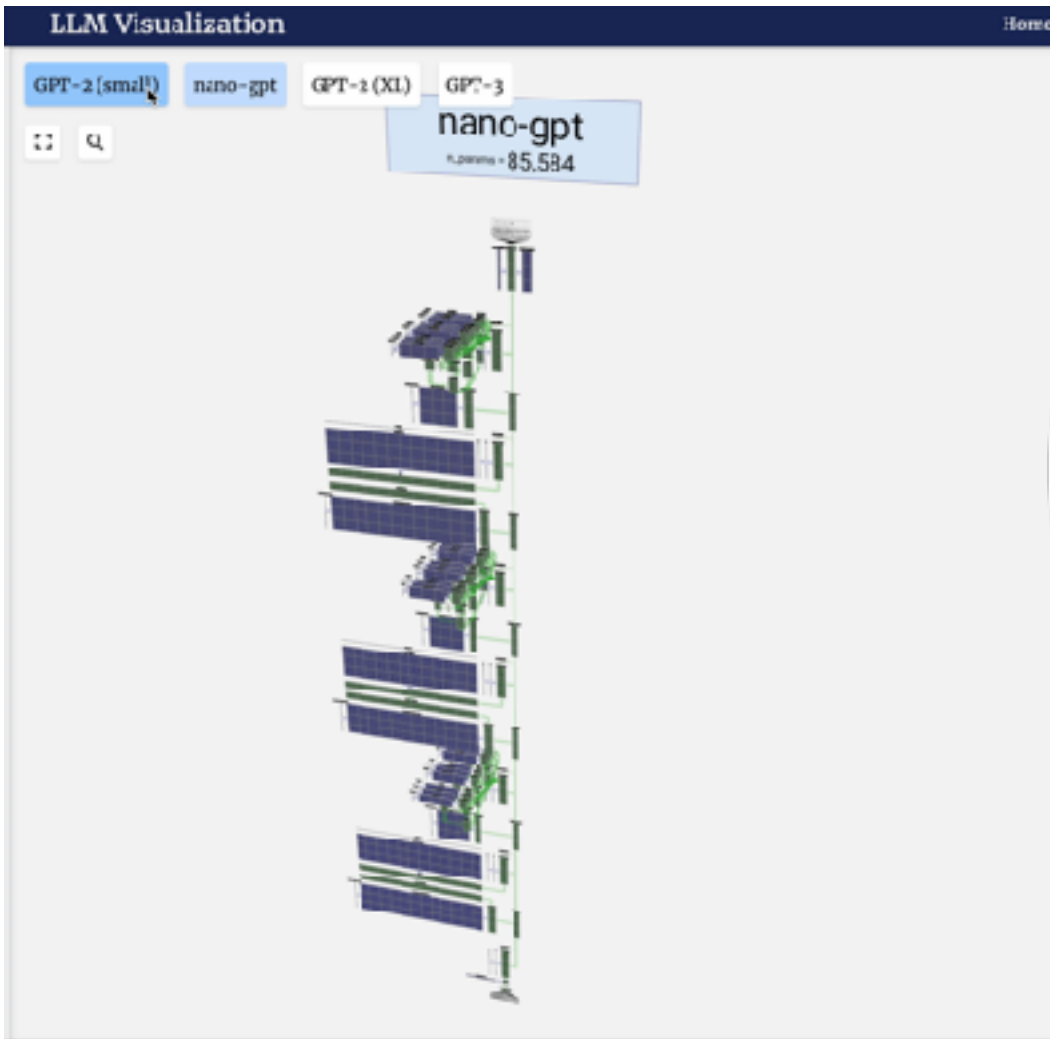
Learn to normalize the rows of Z



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:

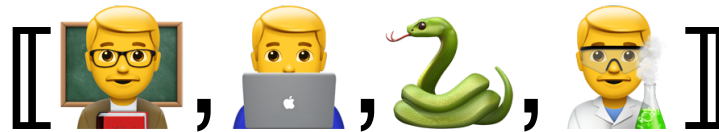




The Transformer and 20
news groups with GloVe

13a. Sequence Basics [Experimental].ipynb

Lecture Notes for **Machine Learning in Python**



Professor Eric Larson
Sequential CNN and Transformers