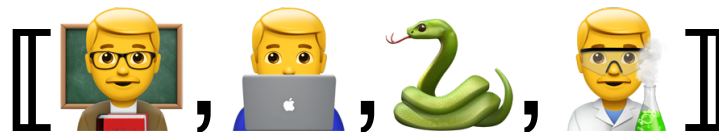


Lecture Notes for **Machine Learning in Python**



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

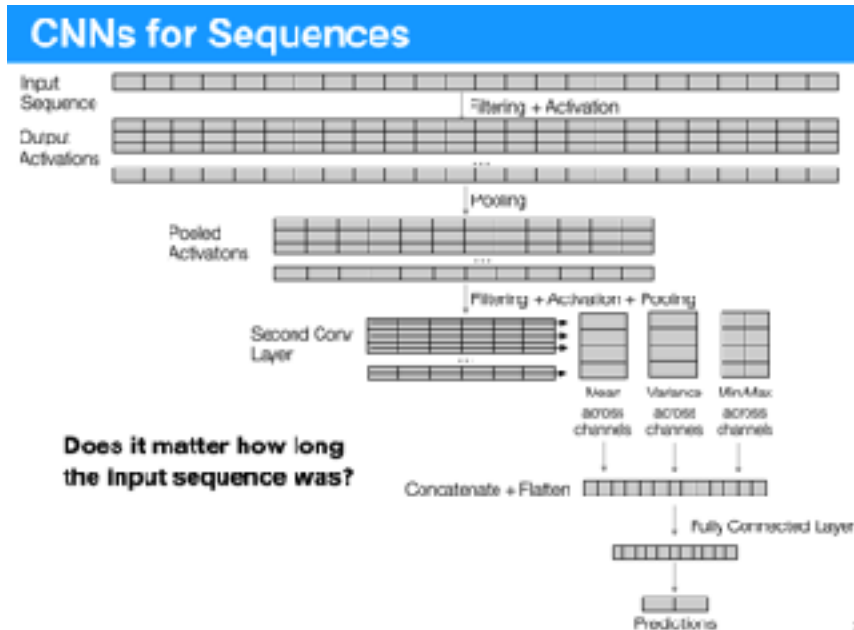
Sequential
Networks

Keras, Tensorflow
Intuition, Detailed implement.

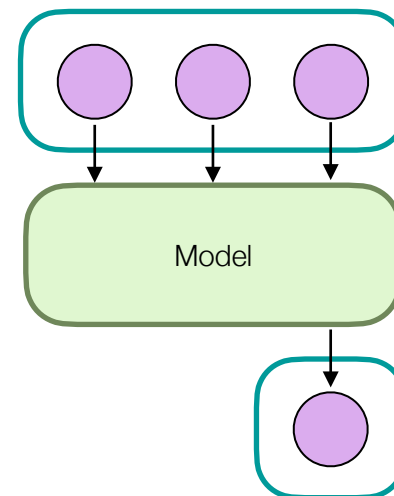
Ethics in
Language Models

ConceptNet
Case studies

Last Time:

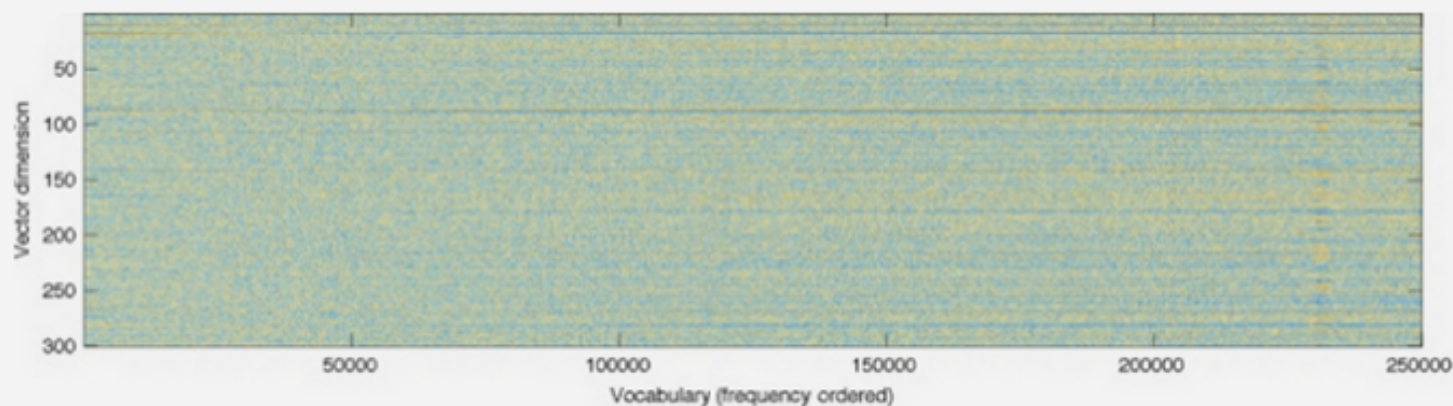


Many to One



Visualization

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

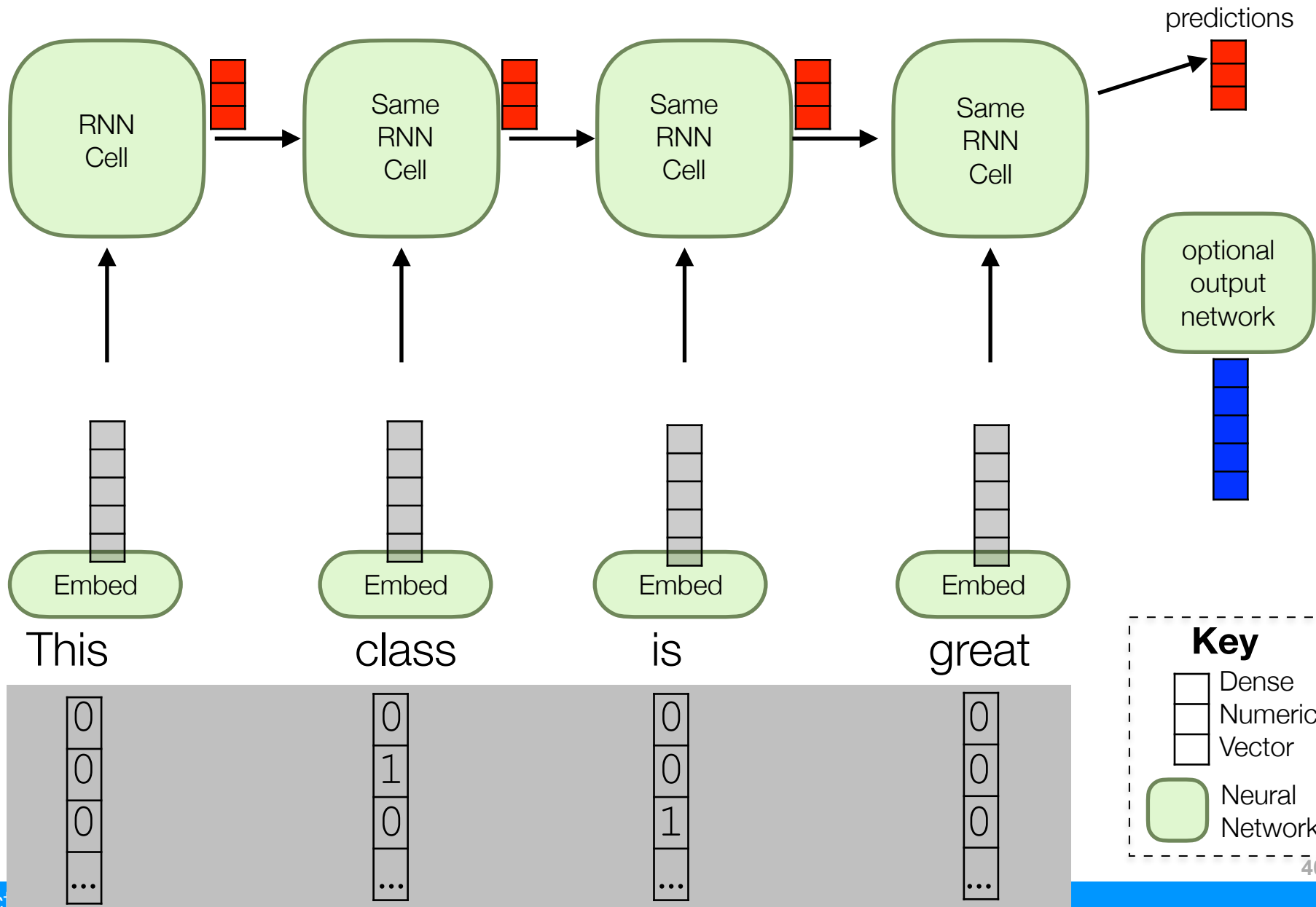


RNNs for Sequences



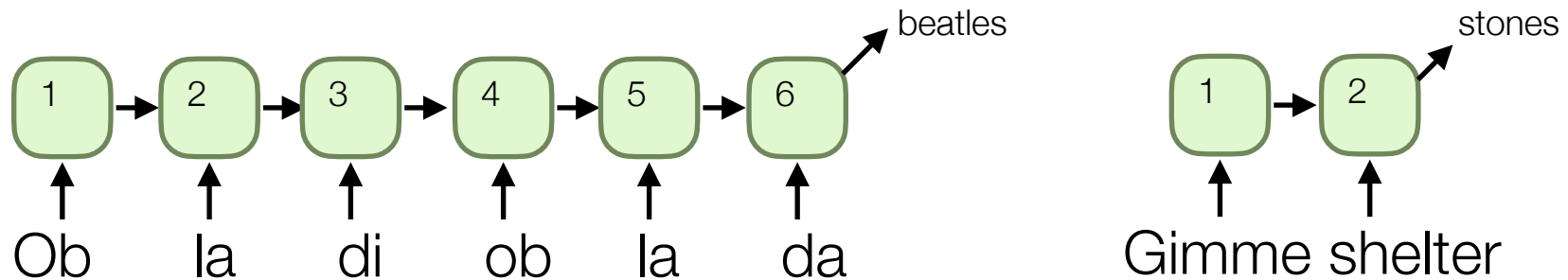
"And how does it make you feel when she jumps over you and calls you a lazy dog?"

Recurrent flow with embeddings



Different length input documents?

- option A: dynamic length sequences



- option B: padding/clipping



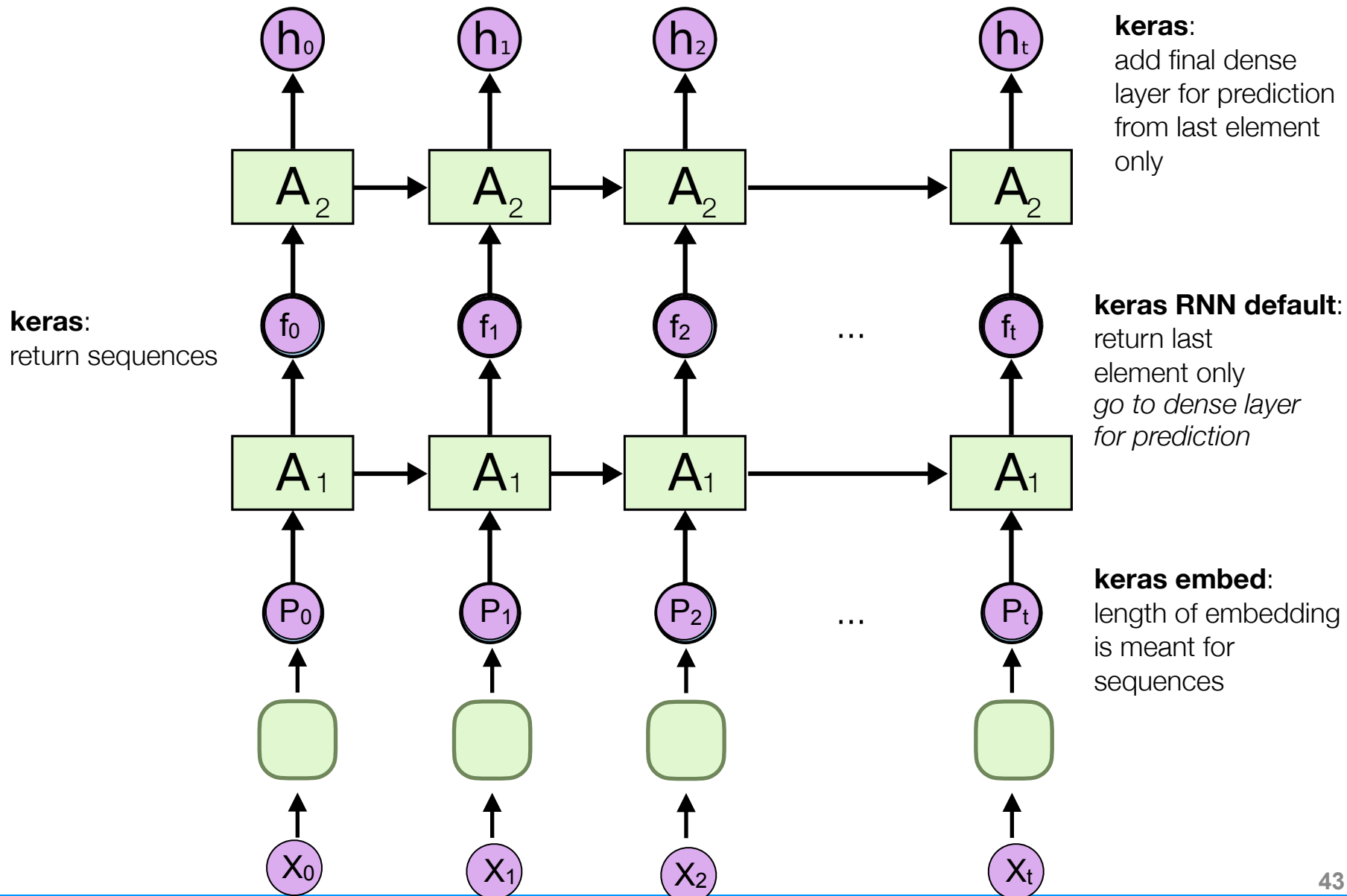
- main difference:

speed based on computation graph design

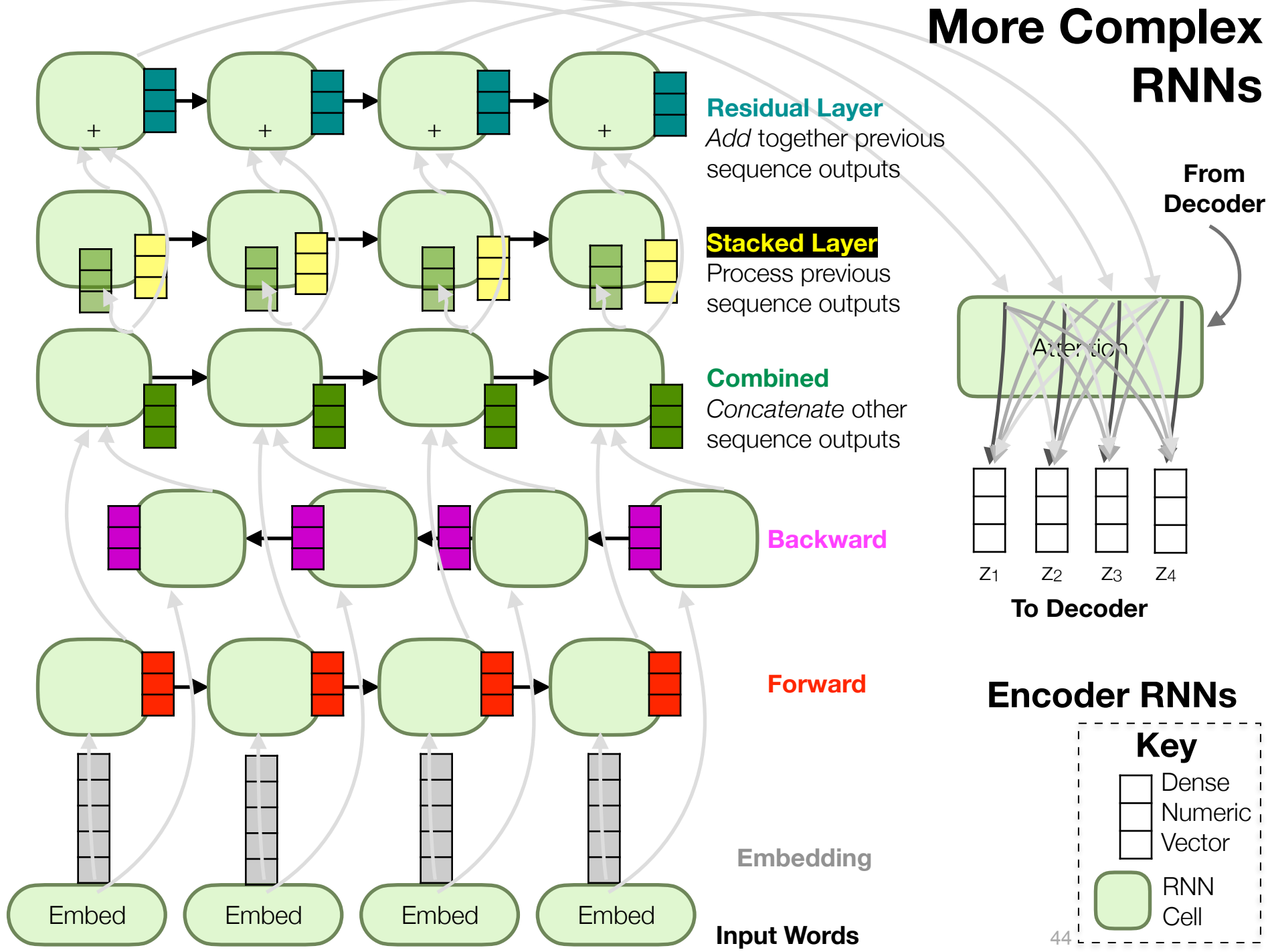
Self Test

- The main reason **dynamic length is slow** is because:
 - A. the computation graph must be updated
 - B. weights must be tied together for each recurrent/sequential node differently
 - C. the embedding matrix cannot be applied in parallel to each word
 - D. no reason: dynamic length is actually faster

Sequence Stacking

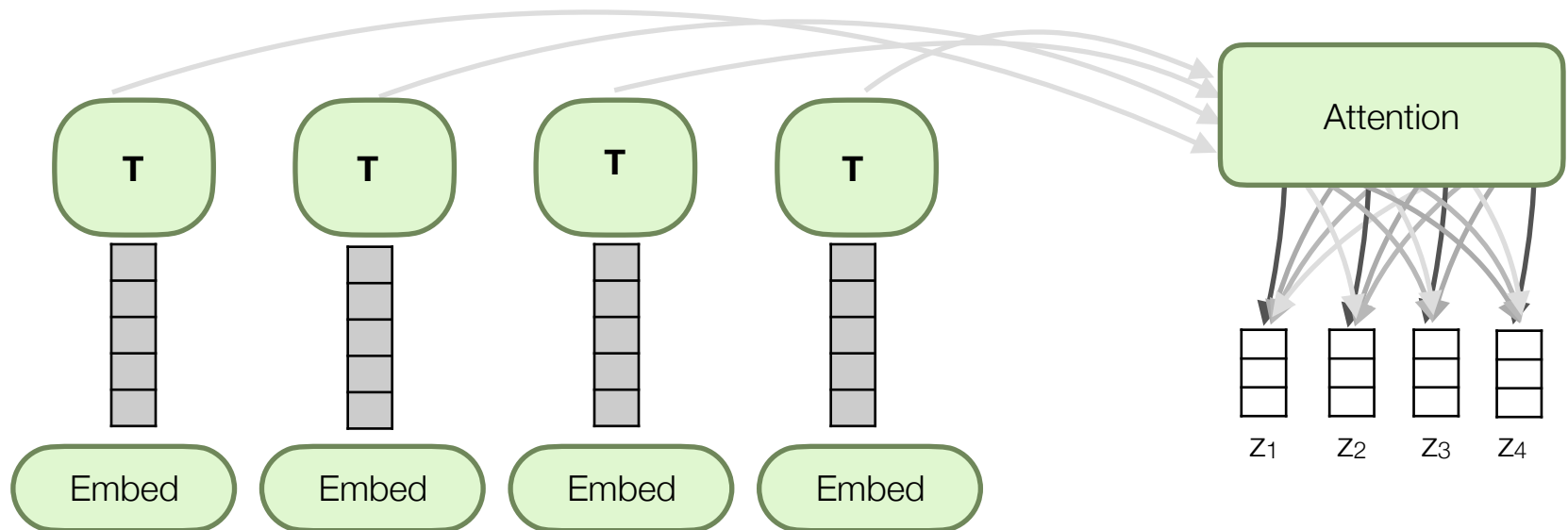


More Complex RNNs



Transformers Intuition (reminder)

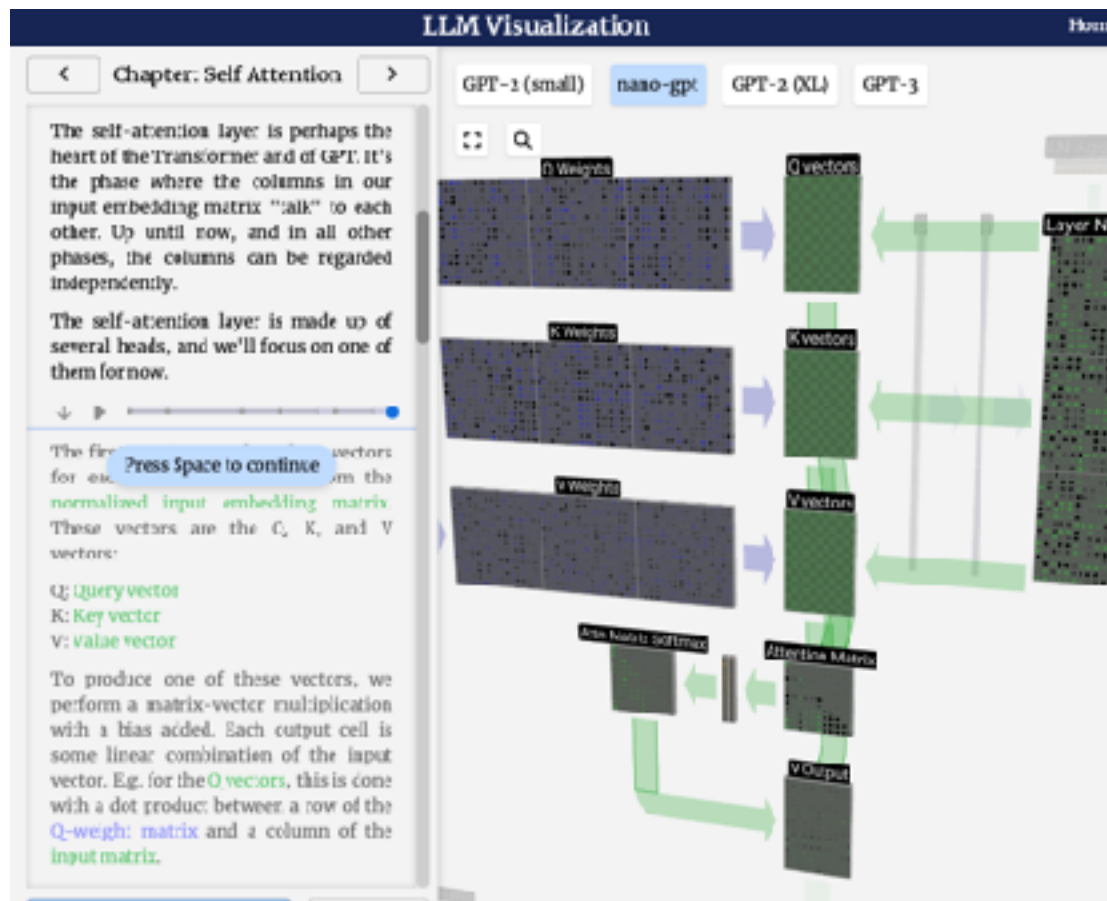
- Recurrent networks track state using an “updatable” state vector, but this takes processing iterative
- Attention mechanism (in RNNs) already takes a weighted sum of state vectors to generate new token in a decoder
- ... so why not just use attention on a transformation of the embedding vectors? **Do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



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

<https://bbycroft.net/llm>

Transformers



Attention is All You Need

- **Continued Motivation:**

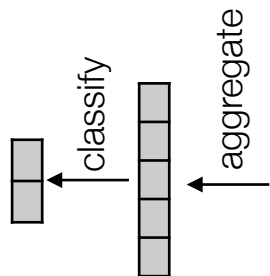
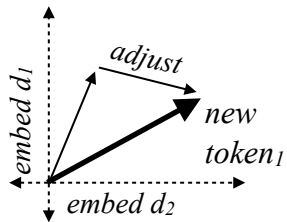
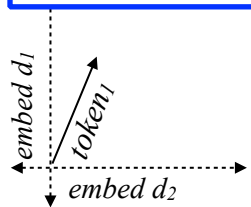
- RNNs are not inherently parallelized or efficient at remembering based on state vector
- CNNs are not resilient to long-term word relationships, limited by filter size

- **Transformer Solution:**

- Build attention into model from the **beginning**
- Compare all words to each other through **self-attention**
- Only update **existing embedding space**, via residual
- Define a notion of “**position**” in the sequence
- ***Should be resilient to long term relationships and be highly parallelized for GPU computing!!***

Transformer Overview

Geometric Overview



Matrix Overview

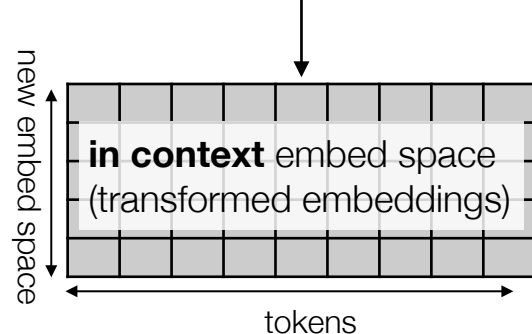
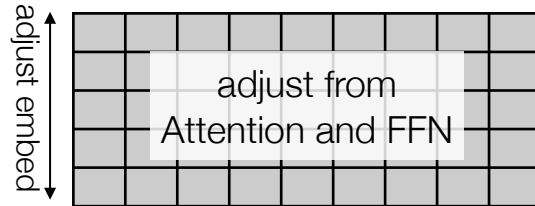
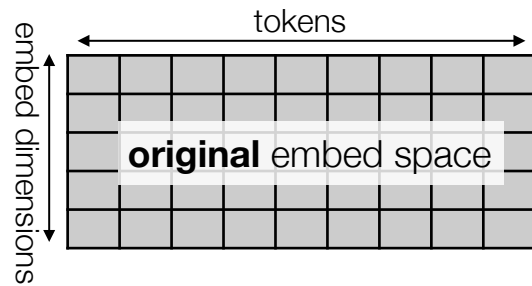
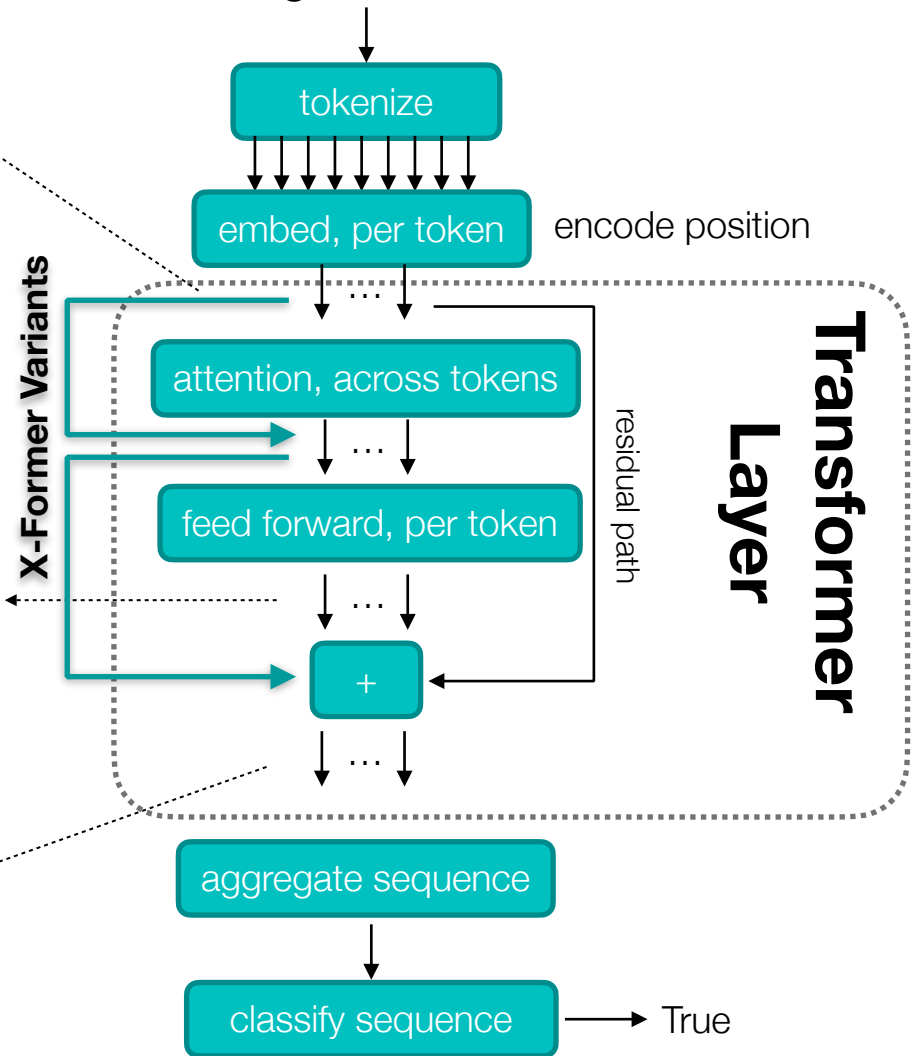


Diagram Overview

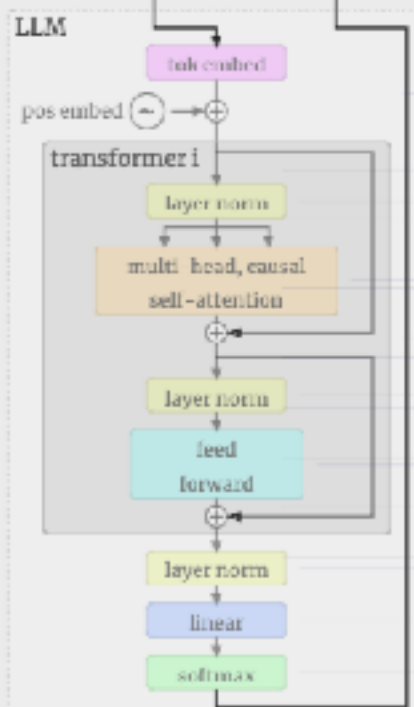
I am waiting for a third AI winter



Auto Regressive Transformer

<https://bbycroft.net/llm>

How to predict text
 14,37 284 4,133
 tokens 14,376
 words 24,56



Components

Embedding
 Layer Norm
 Self Attention
 Projection
 MLP
 Transformer
 Softmax
 Output

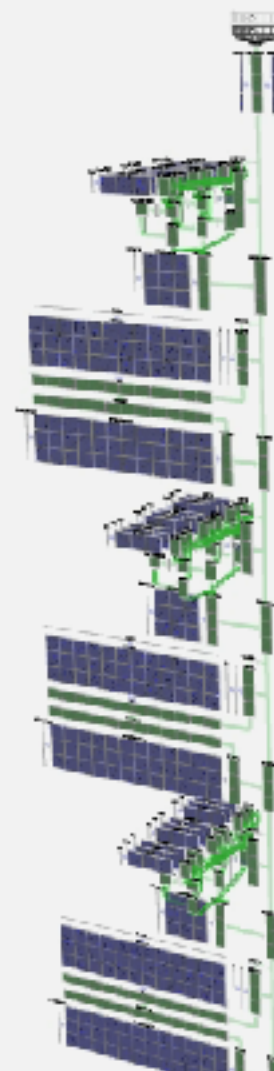
GPT-2 (small)

nano-gpt

GPT-2 (XL)

GPT-3

parameters = 85,584



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

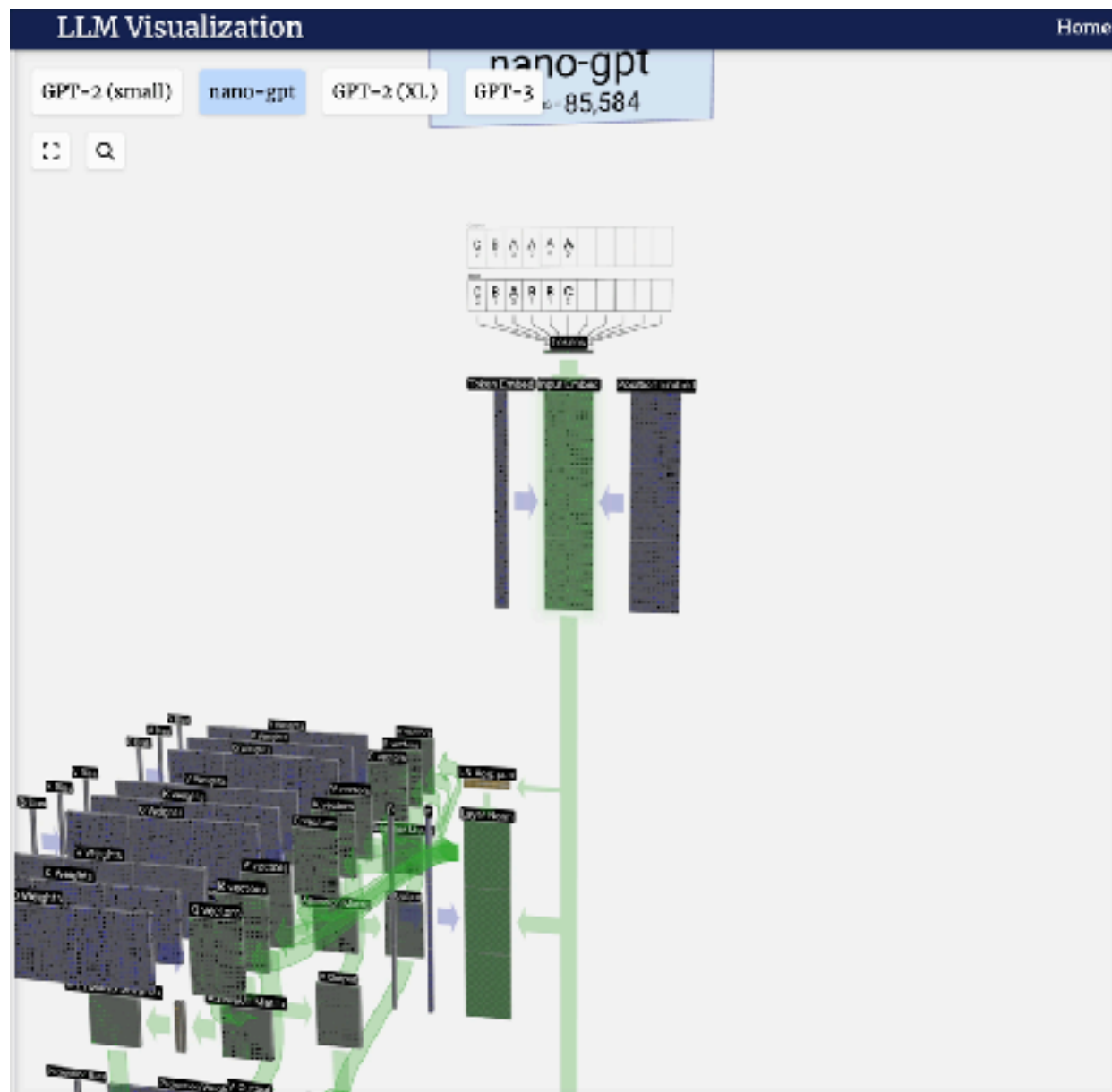
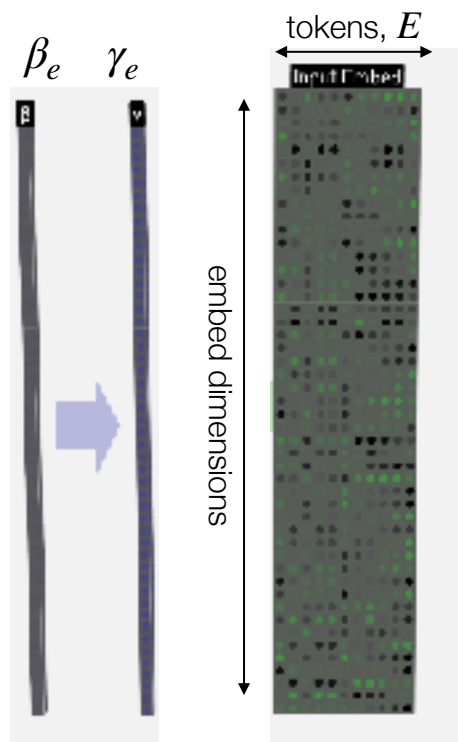
$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

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

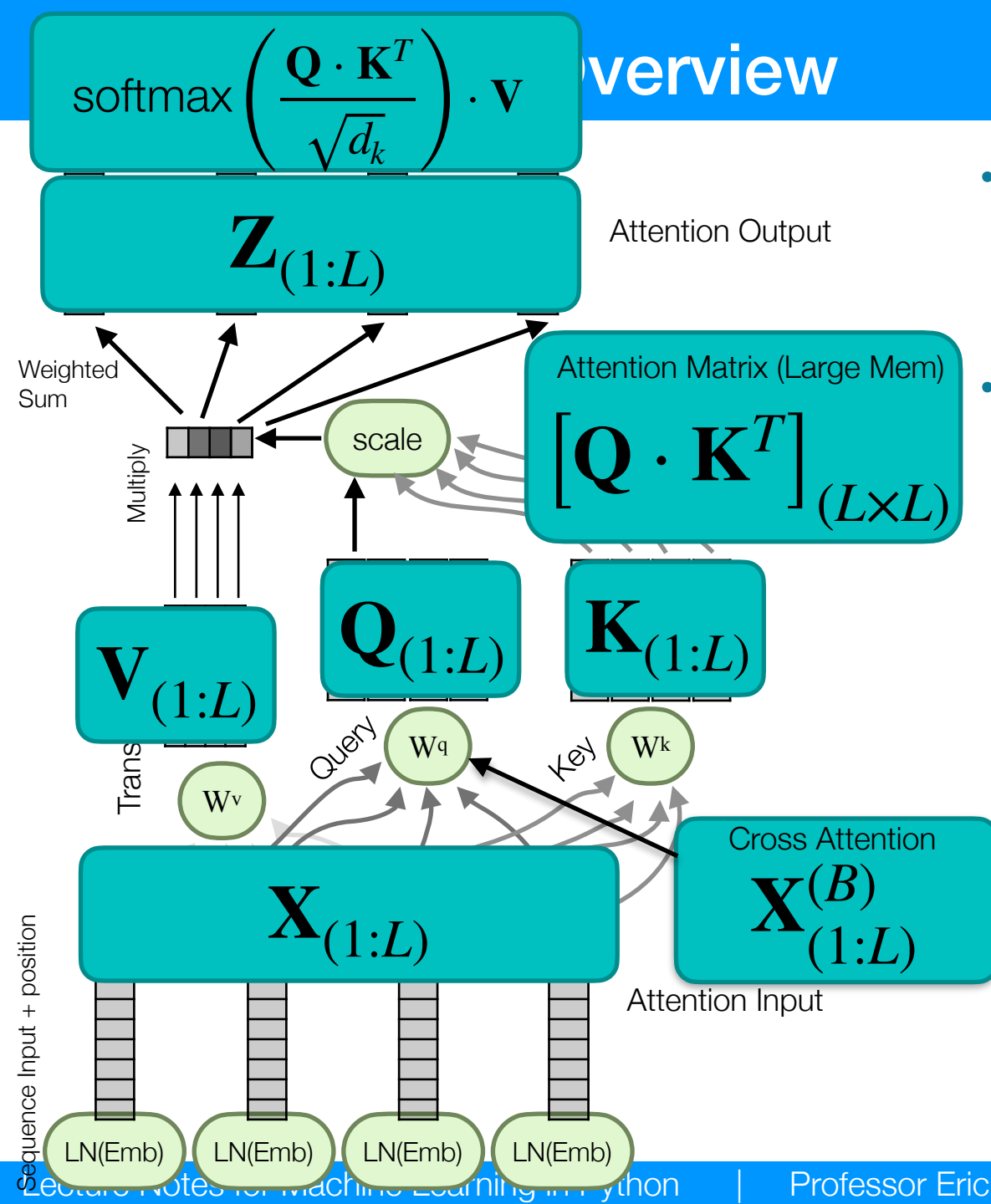
Layer Norm

$$LN(E_{tok}) = \gamma_e \left(\frac{E_{tok} - \mu_{tok}}{\sqrt{\sigma_{tok}^2 + \epsilon}} \right) + \beta_e$$

\uparrow \uparrow
 z-score per token adjust each embed dimension



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

Self Attention: in more detail

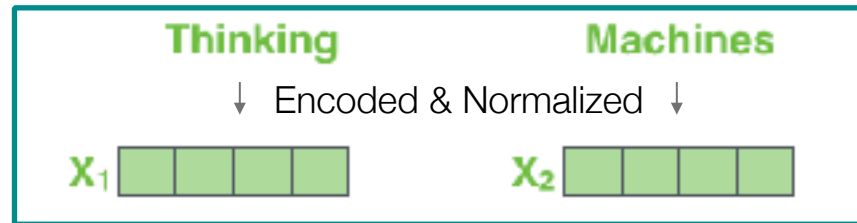
Input

Embedding

Queries

Keys

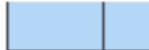

Values



Outputs of Multiplications, per token:

$$q_1 \text{  } = X_1 \cdot W^Q \quad q_2 \text{  } = X_2 \cdot W^Q$$

$$k_1 \text{  } = X_1 \cdot W^K \quad k_2 \text{  } = X_2 \cdot W^K$$

$$v_1 \text{  } = X_1 \cdot W^V \quad v_2 \text{  } = X_2 \cdot W^V$$

Multiply These
(in parallel, for each token)

Learned Matrices



W^Q



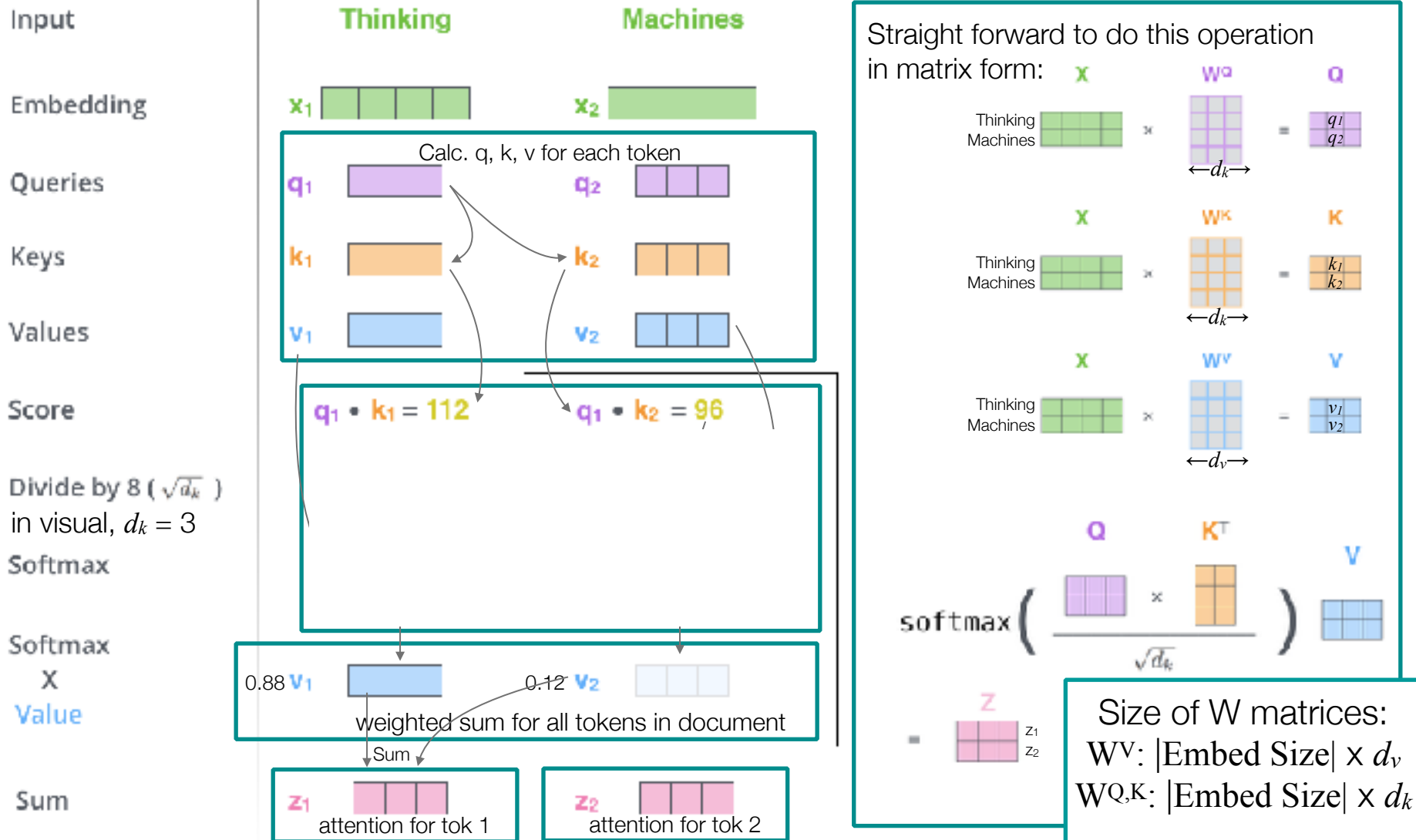
W^K



W^V

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

Self Attention: 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:

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

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

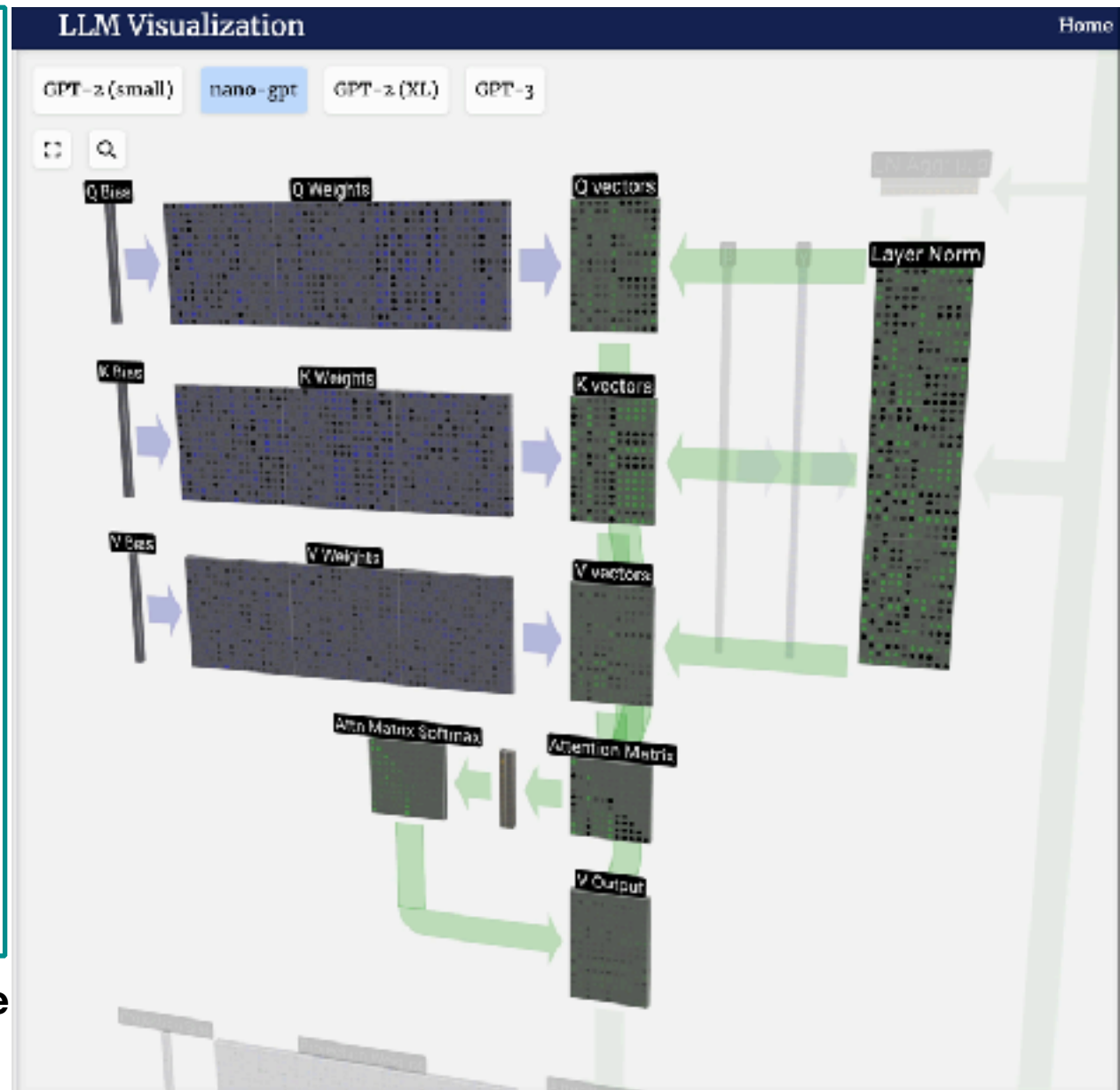
$$\begin{matrix} \text{Thinking} \\ \text{Machines} \end{matrix} \begin{matrix} \mathbf{X} \\ \mathbf{W}^V \end{matrix} = \mathbf{V}$$

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

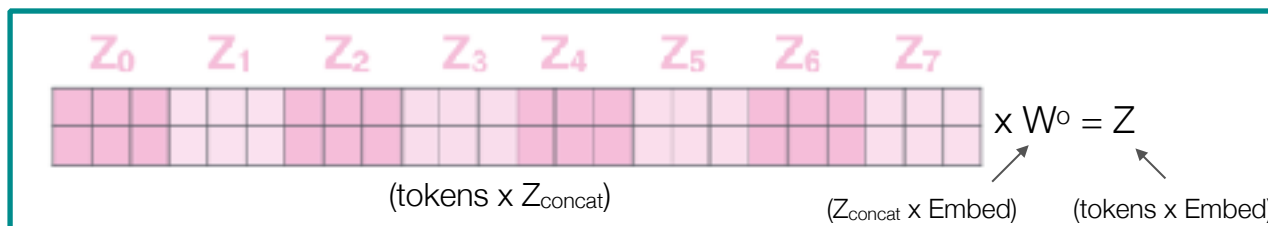
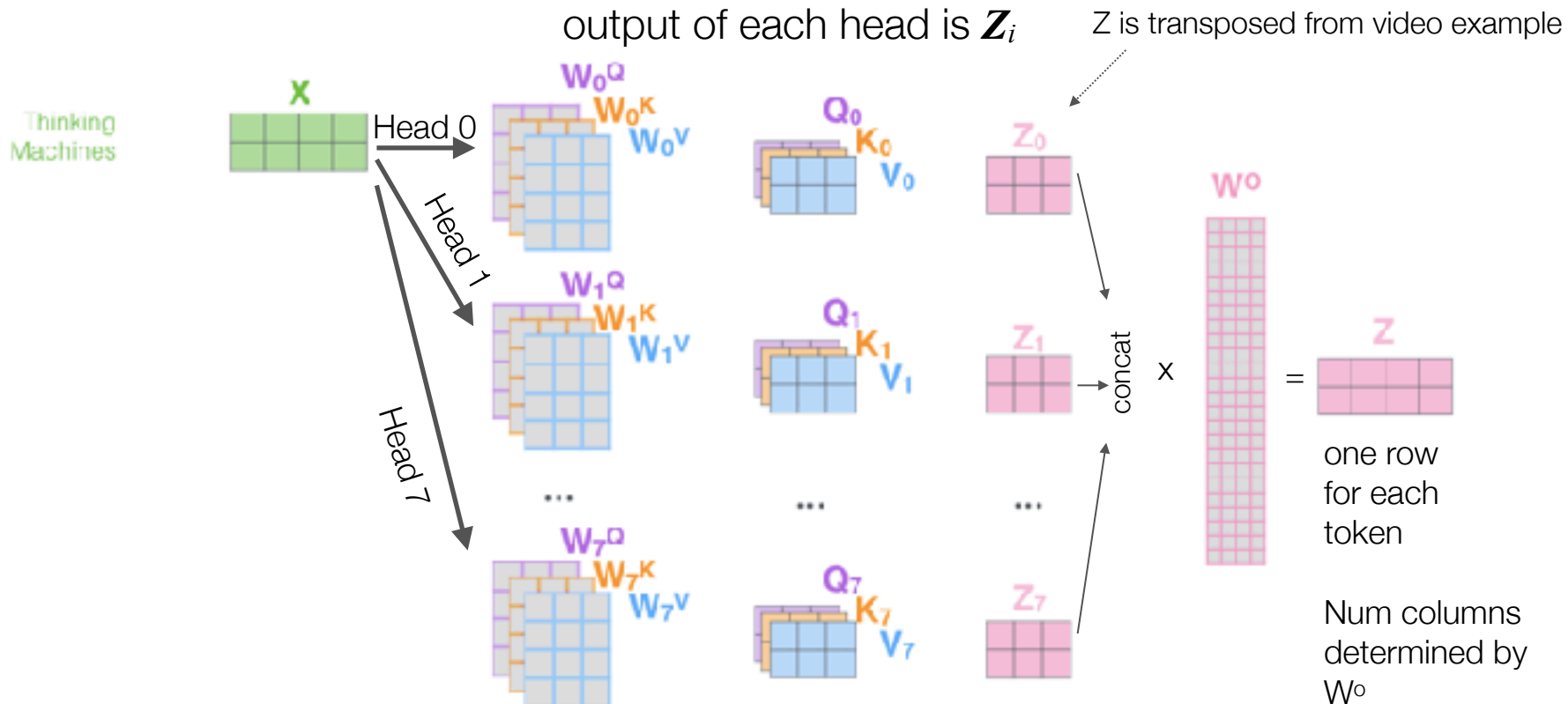
$$\mathbf{Z} = \begin{matrix} \mathbf{Z}_1 \\ \mathbf{Z}_2 \end{matrix}$$

output of each head is \mathbf{Z}_i

After Self-attention, **heads are processed by a FFN**

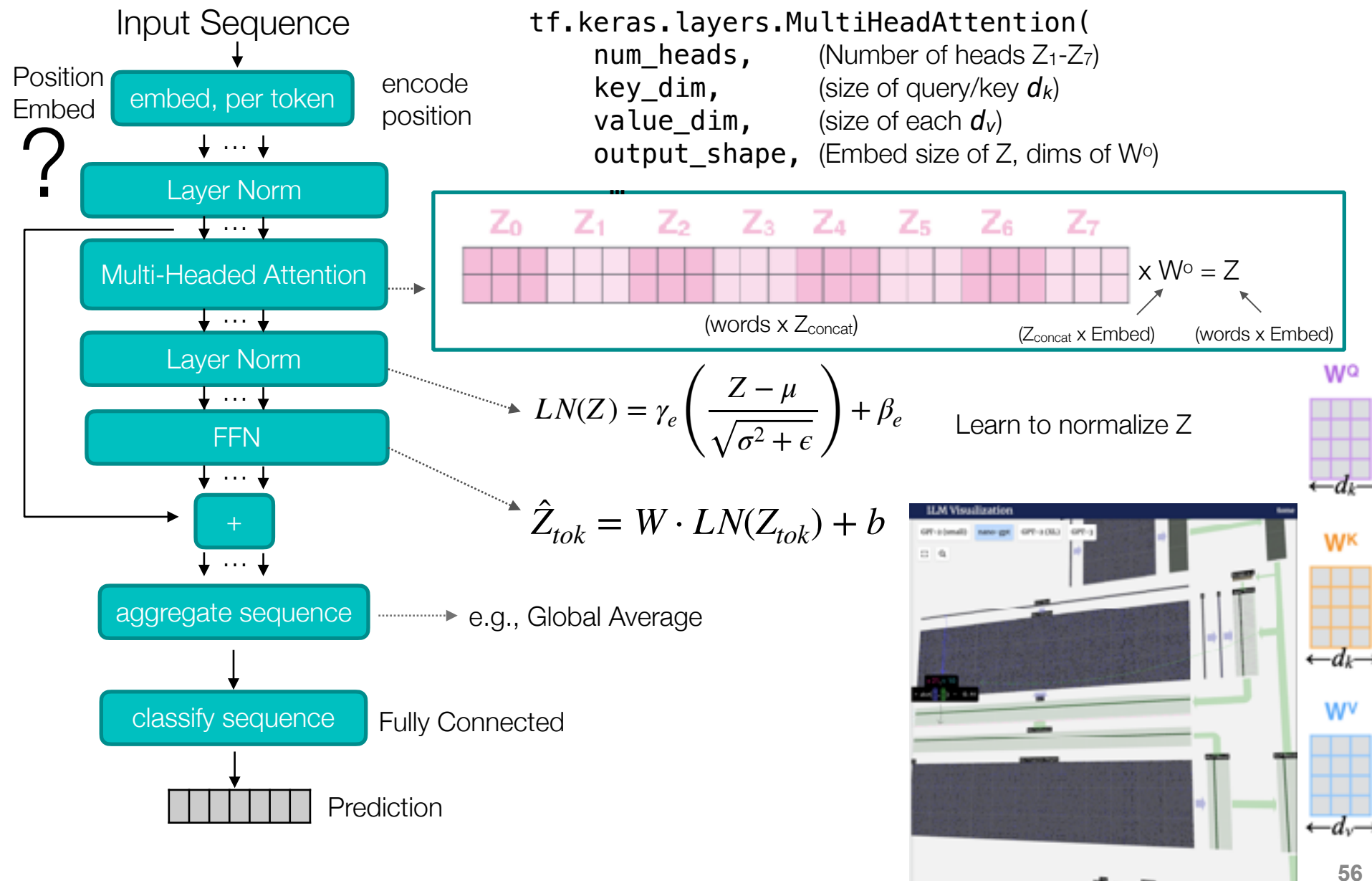


Transformer: Multi-headed Attention



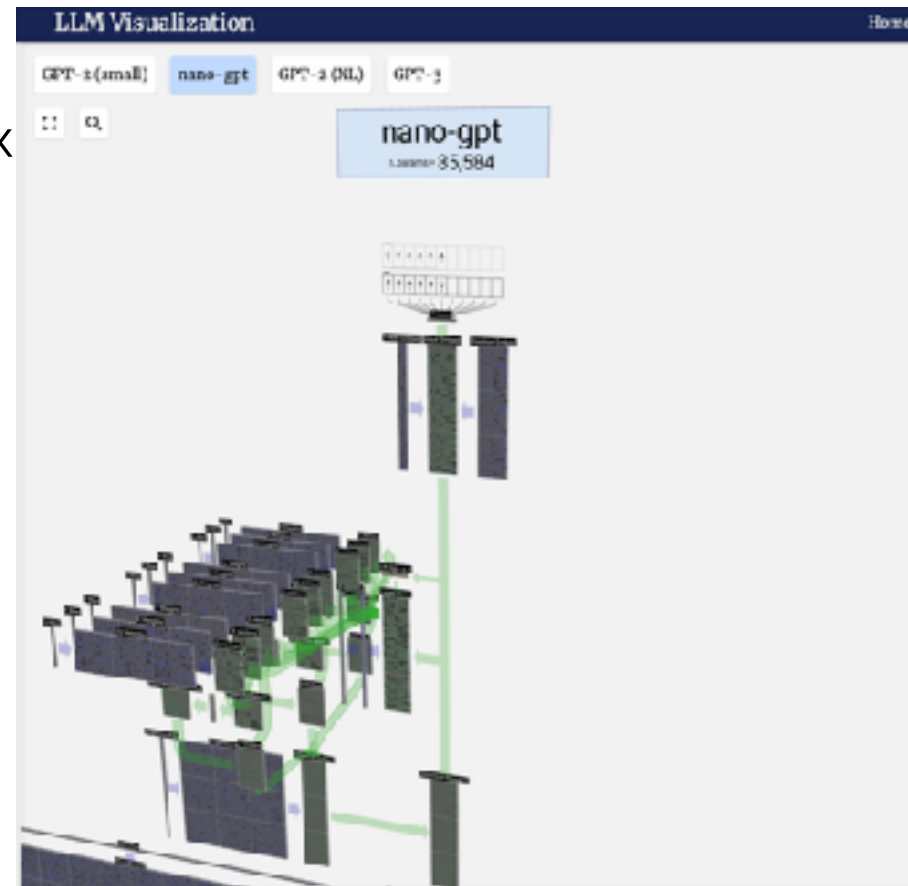
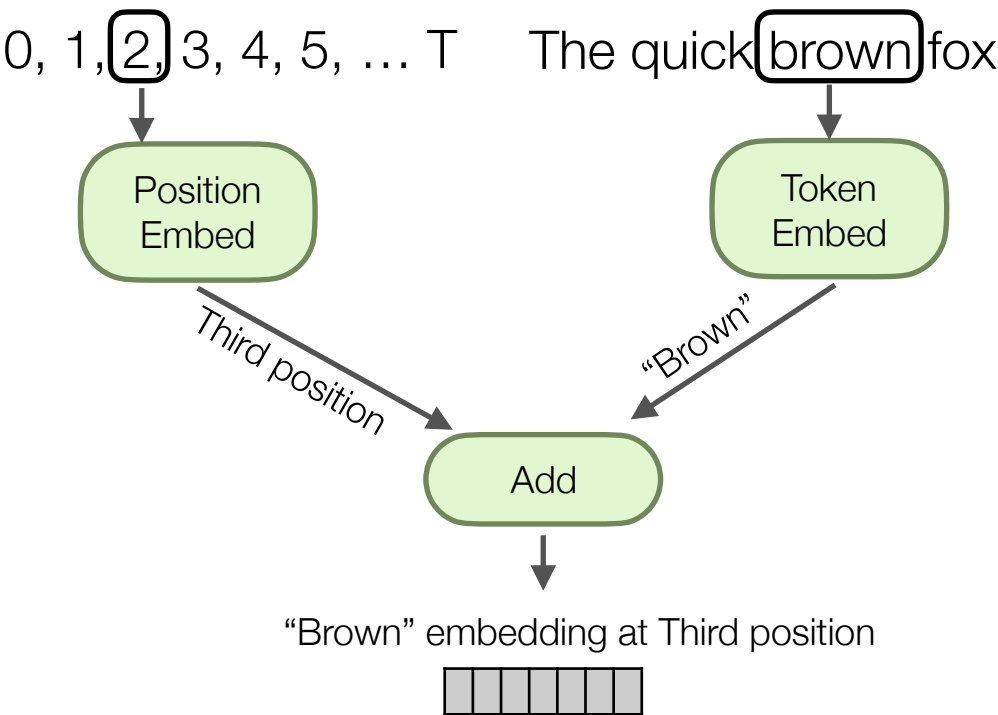
Excellent Blog on Transformers: <http://jalammar.github.io/illustrated-transformer/>

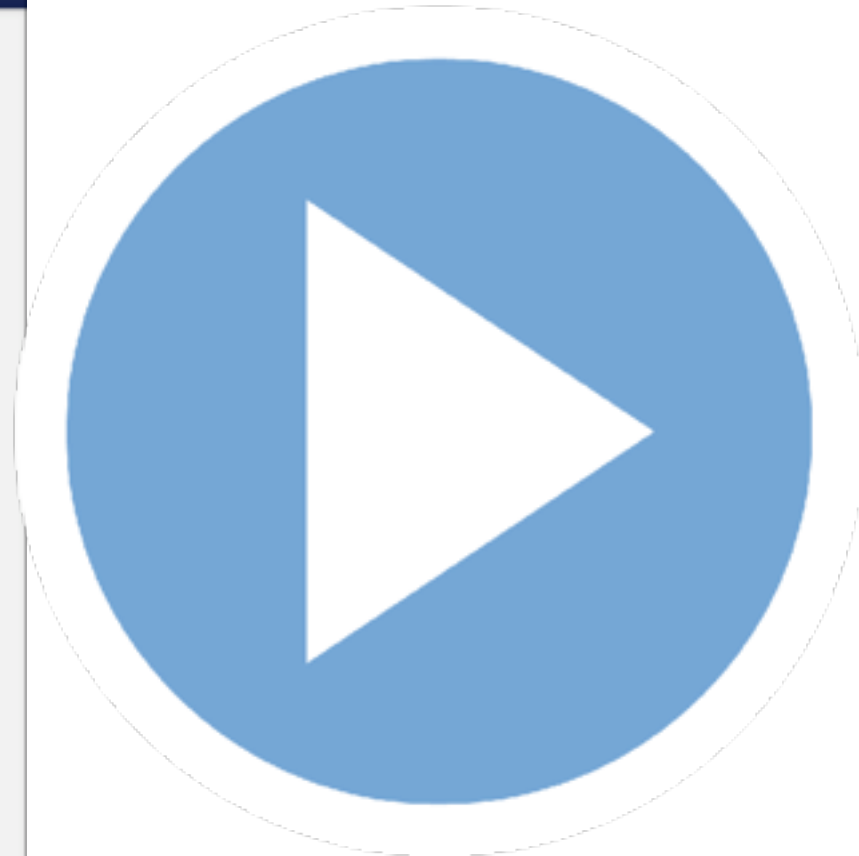
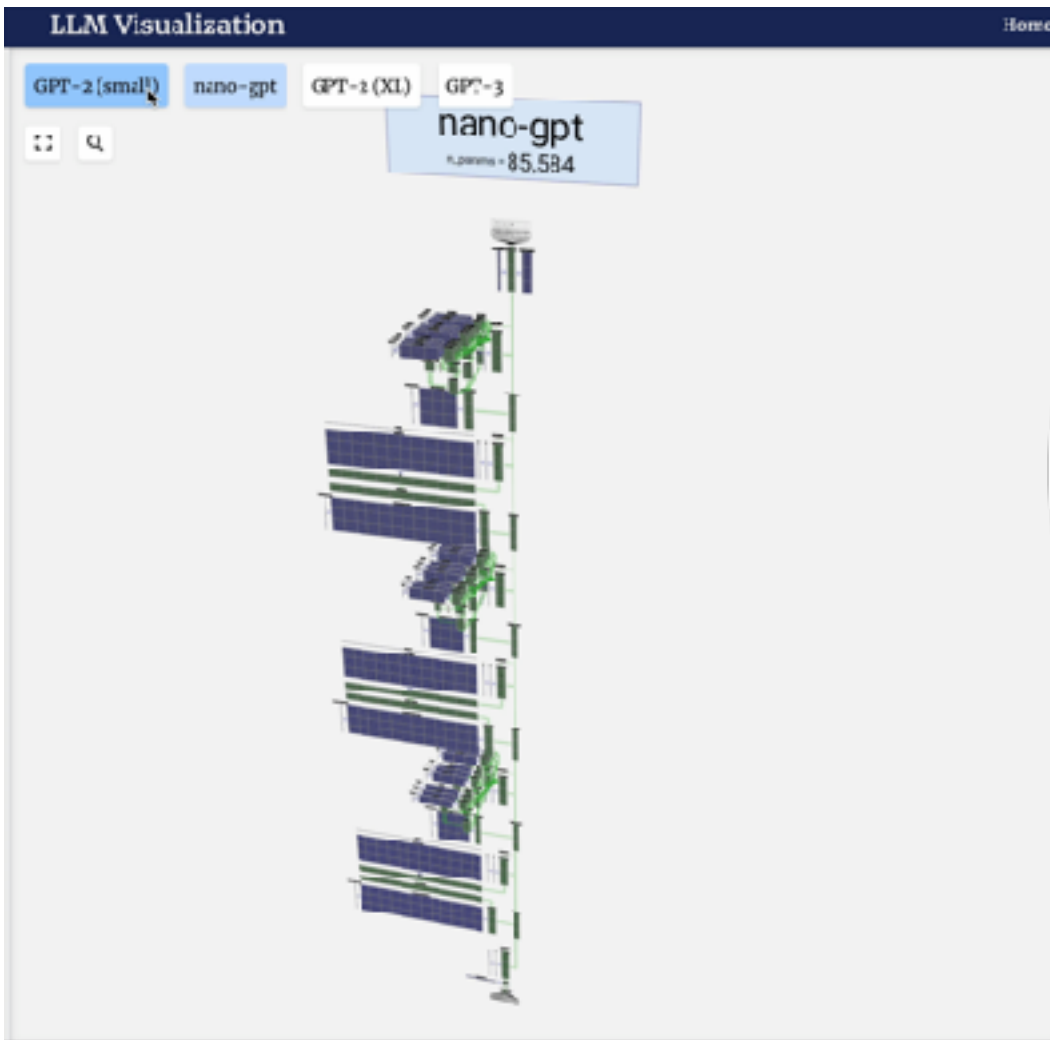
Putting It Together



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:

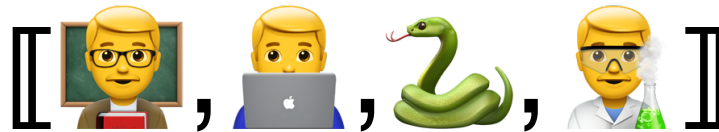




The Transformer and 20
news groups with GloVe

13a. Sequence Basics [Experimental].ipynb

Lecture Notes for **Machine Learning in Python**



Professor Eric Larson
Sequential CNN and Transformers