Lecture Notes for
**Machine Learning in Python**

Professor Eric Larson
**Transformer Demo + Positional Encoding**

# Lecture Agenda

- Logistics
    - Grading Update
    - Sequential Networks due **see canvas**
- Agenda
    - Sequential Networks Demo
        - Extended Demo
    - Final Town Hall
    - (if time) More effective position encoding
    - (if time) More efficient attention
    - Next time (if not behind):
        - Finish above lecture topics
        - Retrospective and Evaluations

61

| 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
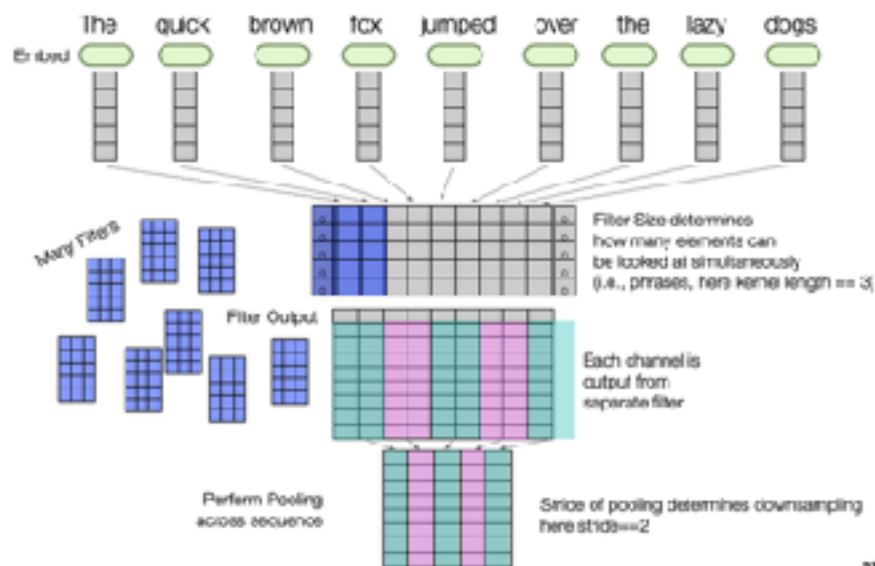Detailed mathematics for NN operations

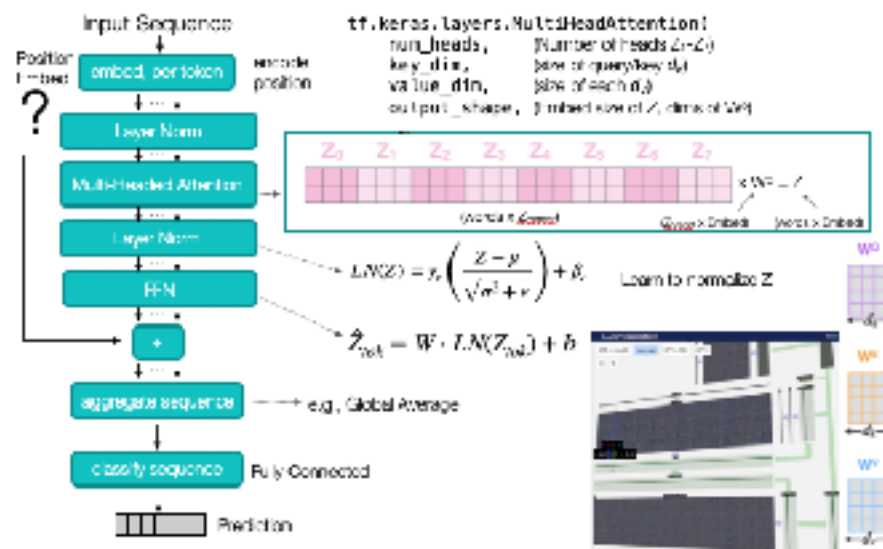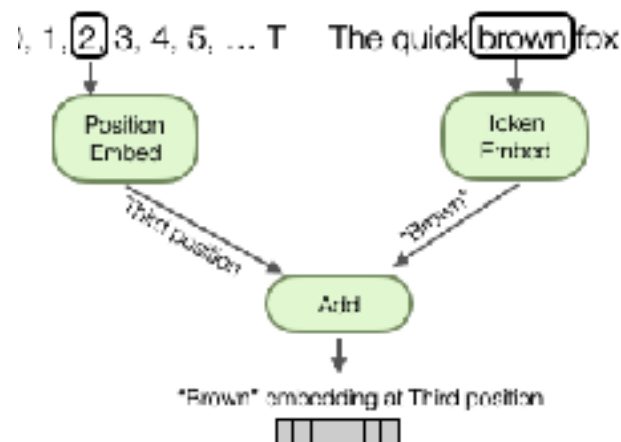| Wide and Deep Networks | Convolutional Networks | Sequential Networks |
|---|---|---|

Keras, Tensorflow
Intuition, Detailed implement.

Ethics

Input sequence of words and their positions

tokenize

position

embed, per token

+

embed, per position

...

attention, across tokens

...

feed forward, per token

...

+

...

residual path
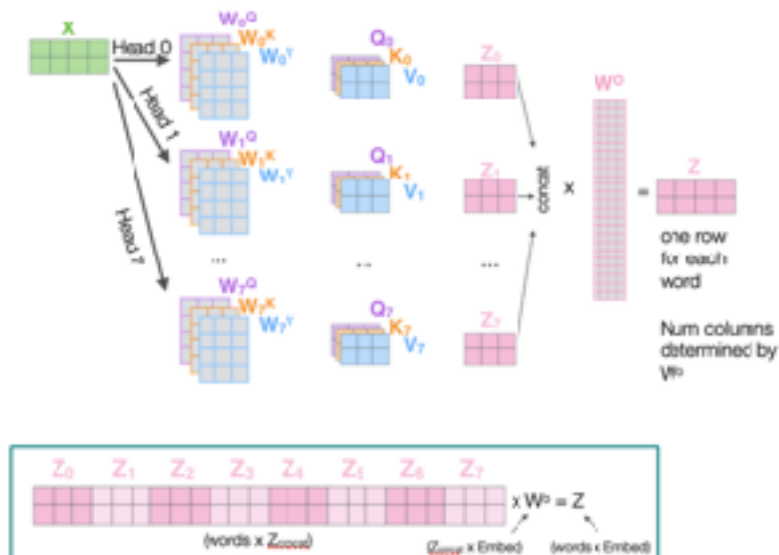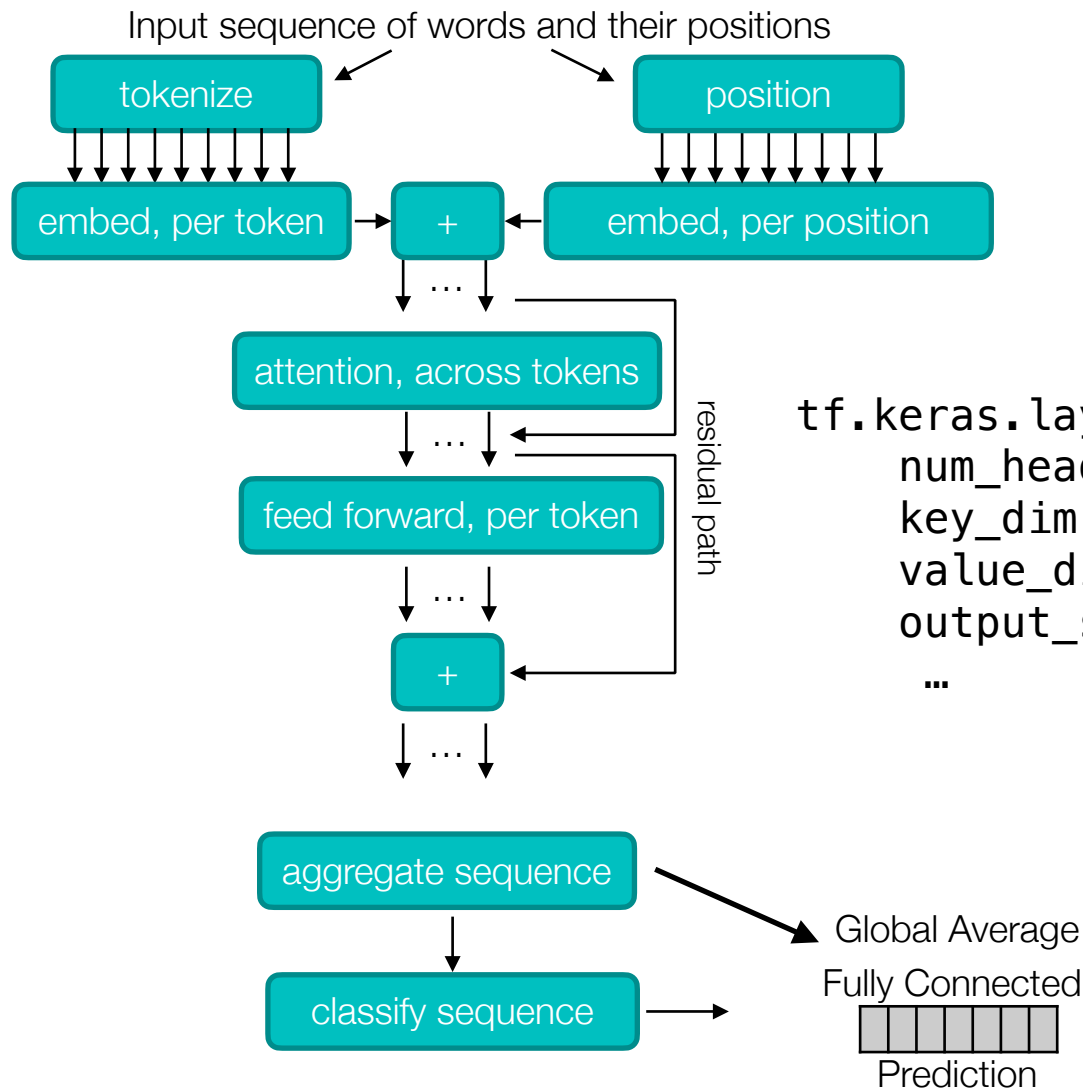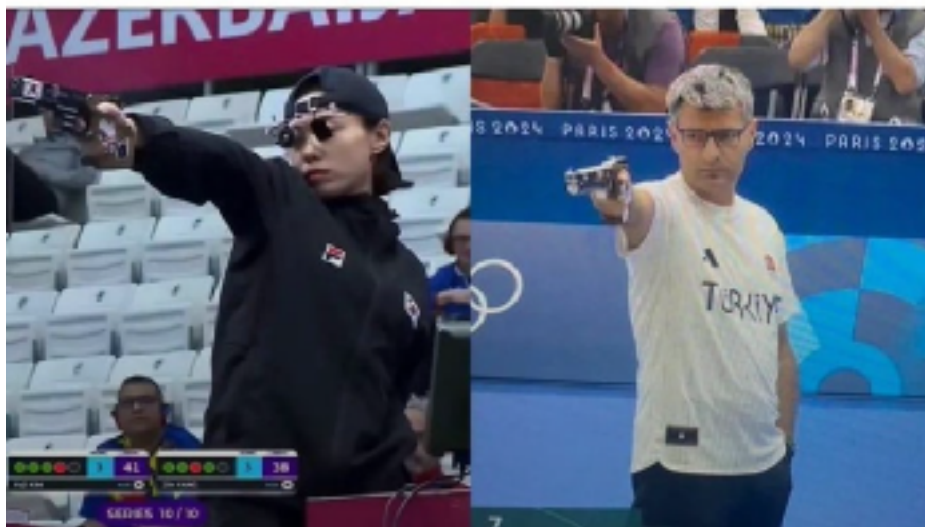
```
tf.keras.layers.MultiHeadAttention(
    num_heads,      (Number of heads $Z_1$-$Z_7$)
    key_dim,        (size of query/key $d_k$)
    value_dim,      (size of each $d_v$)
    output_shape,   (Embed size of Z, dims of $W^o$)
    ...
```

aggregate sequence

Global Average

classify sequence

Fully Connected

Prediction

The Transformer and 20 news groups with GloVe
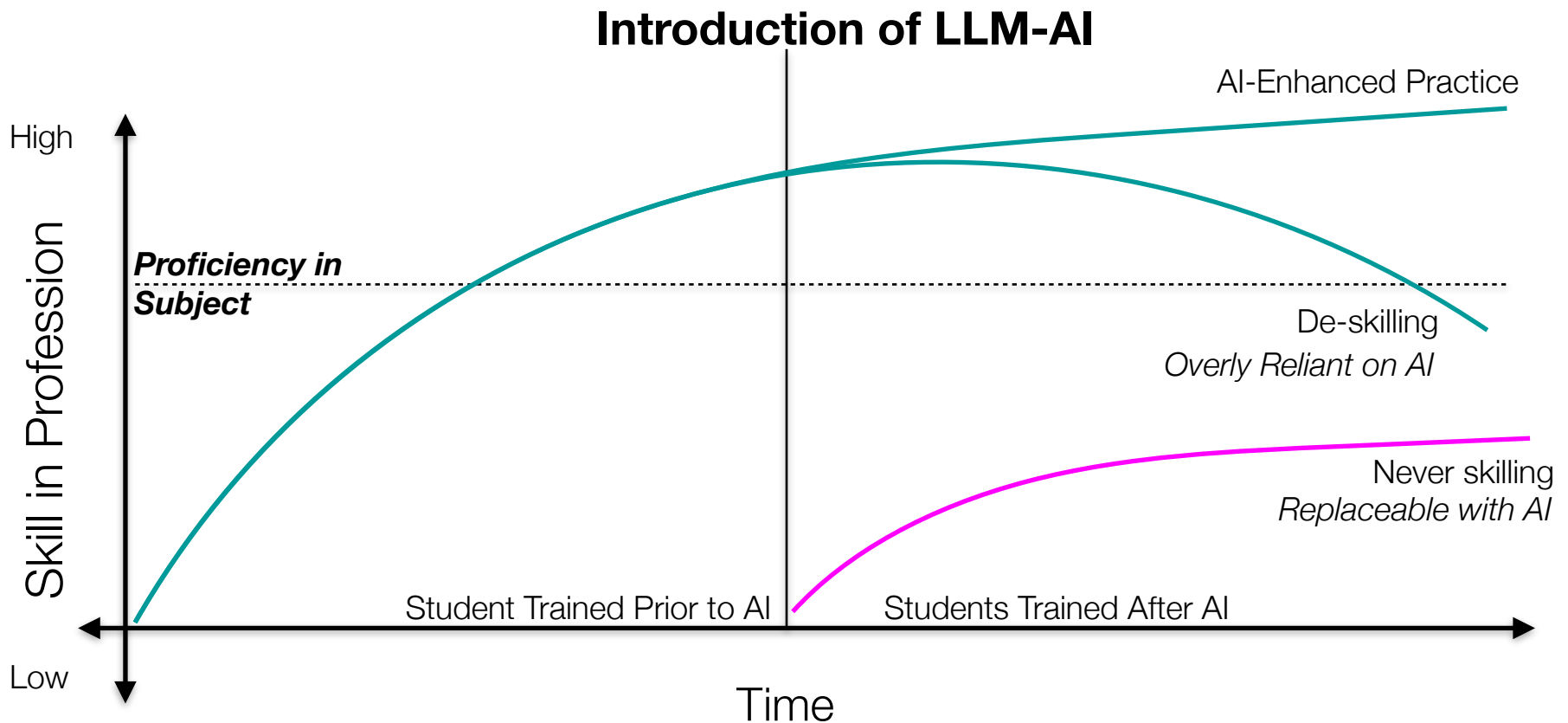
`13a. Sequence Basics [Experimental].ipynb`

# Sequential Networks Town Hall

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

Attention



65

# Positional Encoding



**Introduction of LLM-AI**

AI-Enhanced Practice

High

**Proficiency in Subject**

De-skilling
*Overly Reliant on AI*

Never skilling
*Replaceable with AI*

Skill in Profession

Student Trained Prior to AI

Students Trained After AI

Low

Time

Num Outputs is Same as Inputs

**Do these outputs change, if the input sequence changes order?**

The order of vectors will change, but not the values of each vector…

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

$$\mathbf{V}_{(1:L)} \quad \mathbf{Q}_{(1:L)} \quad \mathbf{K}_{(1:L)}$$

$$\mathbf{X}_{(1:L)}$$

67

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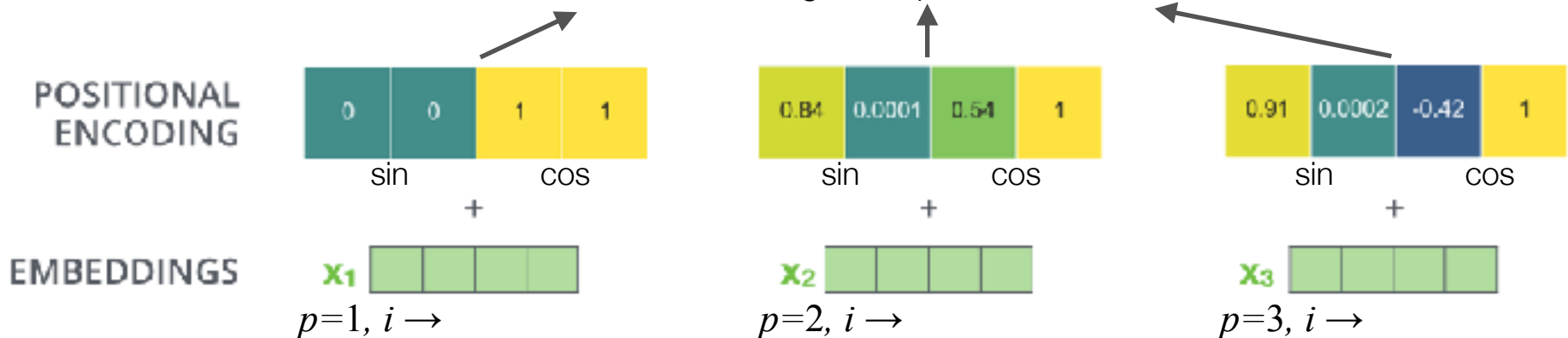- Attempt in paper: add sin/cos to embedding



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

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

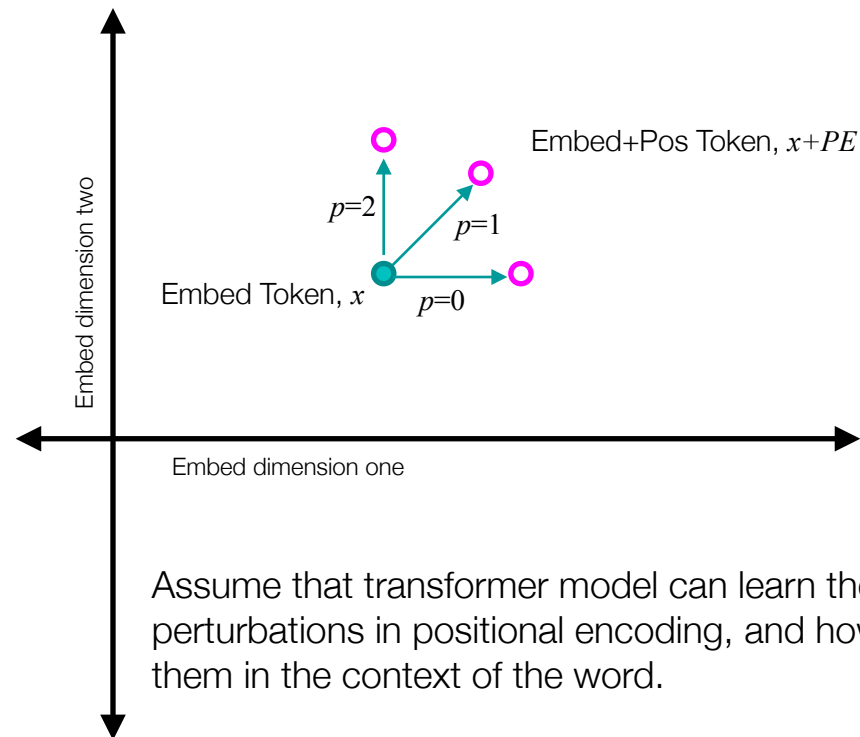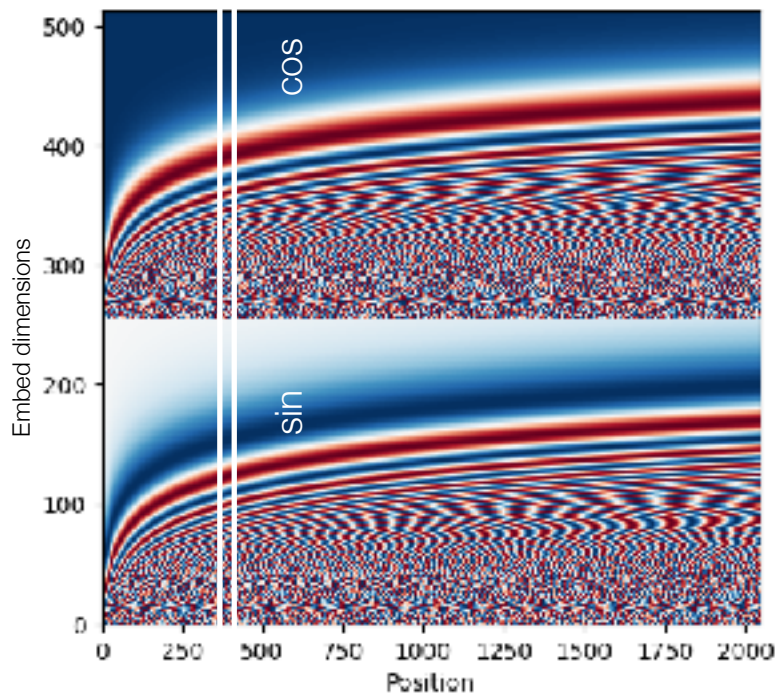$$PE_{(p, i \in d_m...2d_m)} = \cos(p/10000^{(i-d_m)/d_m})$$

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

$X_1$     $p$=1, $i \rightarrow$

$X_2$     $p$=2, $i \rightarrow$

$X_3$     $p$=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.
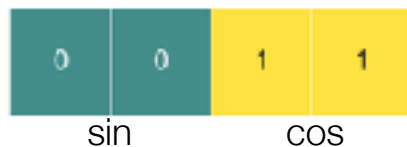
**Excellent Blog on Transformers:** http://jalammar.github.io/illustrated-transformer/

68

# Positional Intuition, Geometrically



Embed+Pos Token, $x+PE$

$p=2$

$p=1$

Embed Token, $x$    $p=0$

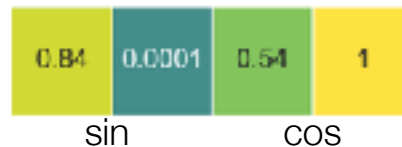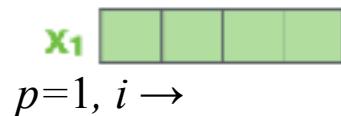Embed dimension two

Embed dimension one

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

POSITIONAL ENCODING

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

+

EMBEDDINGS

$x_1$

$p=1, i \rightarrow$

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

+

$x_2$

$p=2, i \rightarrow$

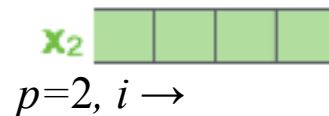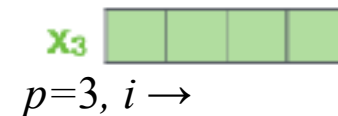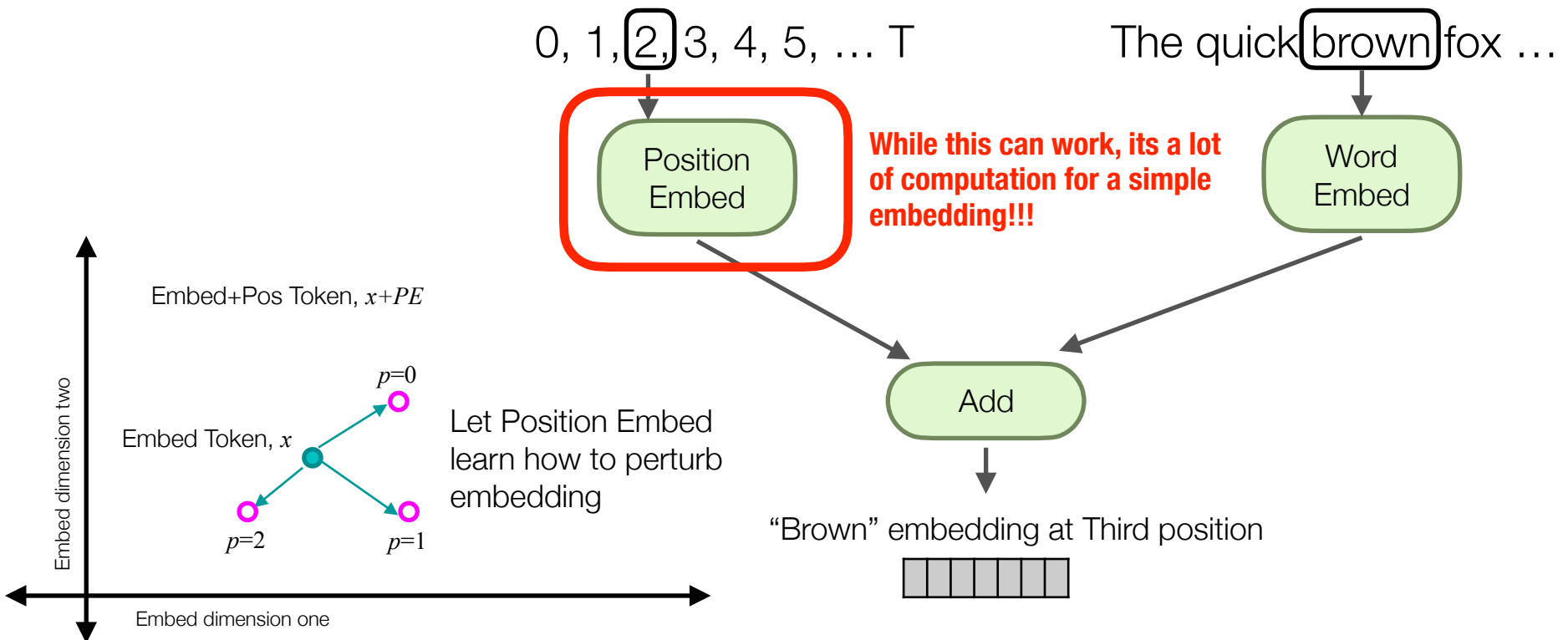| 0.91 | 0.0002 | -0.42 | 1 |
|---|---|---|---|
| sin | | cos | |

+

$x_3$

$p=3, i \rightarrow$

# Transformer: Positional Embedding

- 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

While this can work, its a lot of computation for a simple embedding!!!

Word Embed

Embed+Pos Token, $x+PE$

Embed dimension two

$p=0$

Embed Token, $x$

Let Position Embed learn how to perturb embedding

$p=2$        $p=1$

Embed dimension one

Add

"Brown" embedding at Third position

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

# Relative Positional Encoding

- Relative position encoding:
  add relative words differences into $\mathbf{Q} \cdot \mathbf{K}^T$
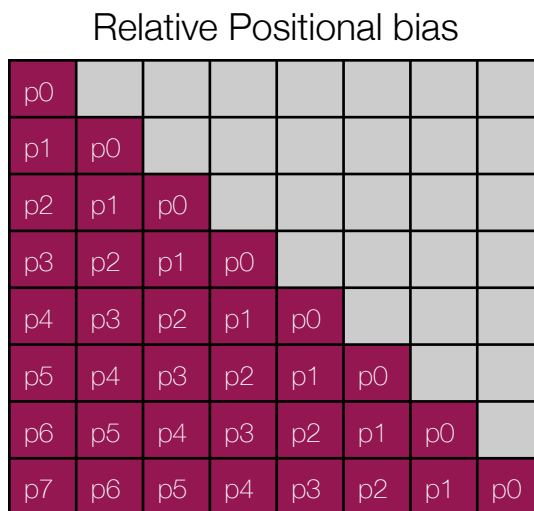
$$\mathbf{Q} \cdot \mathbf{K}^T$$

Relative Positional bias

|    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|
| p0 |    |    |    |    |    |    |    |
| p1 | p0 |    |    |    |    |    |    |
| p2 | p1 | p0 |    |    |    |    |    |
| p3 | p2 | p1 | p0 |    |    |    |    |
| p4 | p3 | p2 | p1 | p0 |    |    |    |
| p5 | p4 | p3 | p2 | p1 | p0 |    |    |
| p6 | p5 | p4 | p3 | p2 | p1 | p0 |    |
| p7 | p6 | p5 | p4 | p3 | p2 | p1 | p0 |

+

- (+) nicely structured position information
- (-) Slow, lots more memory
- (-) fragments ops further, more KV cache misses

- **How might we still encode relative position, without all the overhead?**