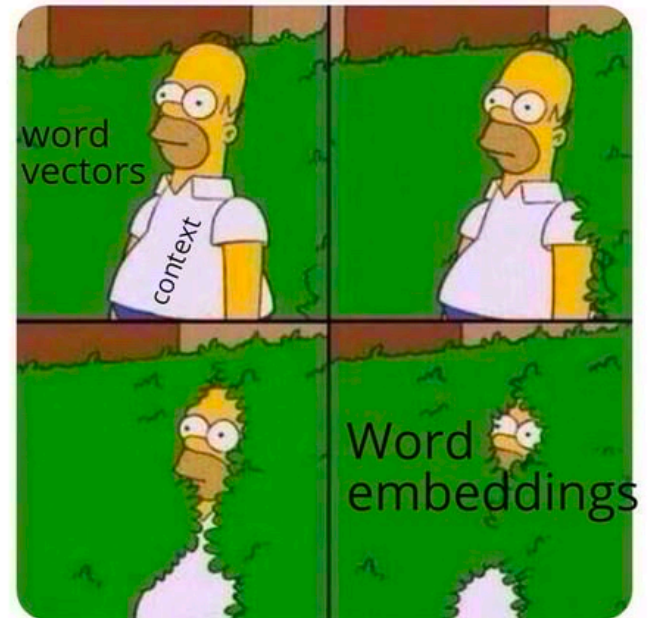May 17, 2023.     DSML : NLP module.     Word embeddings in a nutshell

BERT : Bidirectional Encoder Representation from Transformers.

Class starts @ 9:05

WORDS GET A VECTOR, DOCS GET A VECTOR

EVERYTHING GETS A VECTOR

word vectors

context

Word embeddings

When you penalize your Natural Language Generation model for large sentence lengths

Me think, why waste time say lot word, when few word do trick.

do a little text search and nobody bats an eye

do a little NLP and everybody loses their minds

quickmeme.com

**Recap:**

* **Business Case** : Neural Machine Translation.
→ Given a sentence in English, translate the sentence to French.

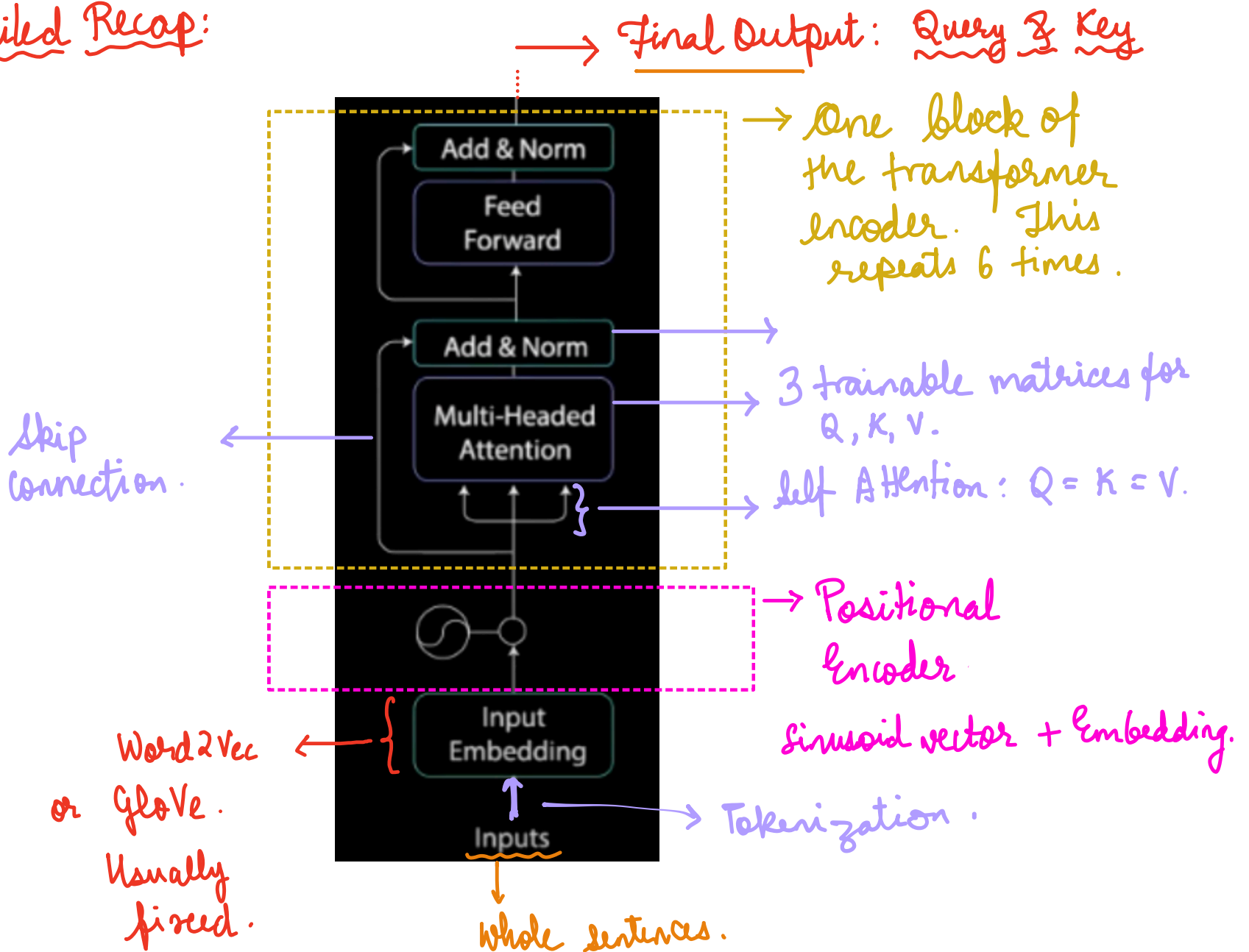* **Approach taken to solve the problem** : Transformer Neural Network.

* **What we covered last class** : Transformer

→ Encoder        Decoder.

* Uses Attention only.
* Requires positional encoding
* Multiheaded self-attention.
* Is trainable.

# Detailed Recap:

→ Final Output: Query & Key
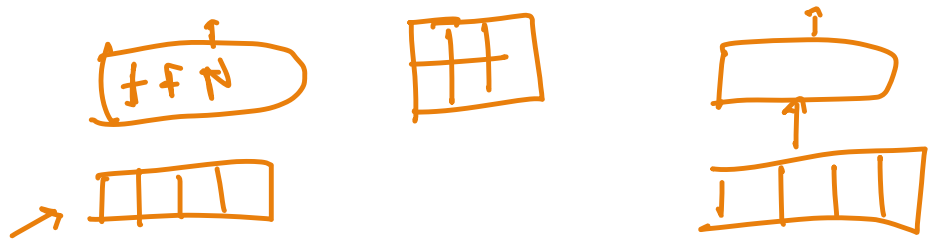


**Add & Norm**

**Feed Forward**

**Add & Norm**

**Multi-Headed Attention**

**Input Embedding**

Inputs

→ One block of the transformer encoder. This repeats 6 times.

Skip connection.

3 trainable matrices for Q, K, V.

Self Attention: Q = K = V.

→ Positional Encoder

Sinusoid vector + Embedding.

Word 2 Vec or GloVe. Usually fixed.

→ Tokenization.

Whole sentences.

# Agenda:

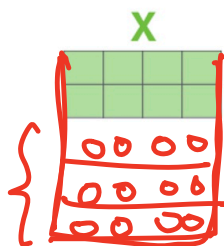✓ * Decoder Architecture.

* BERT: An encoder only model.

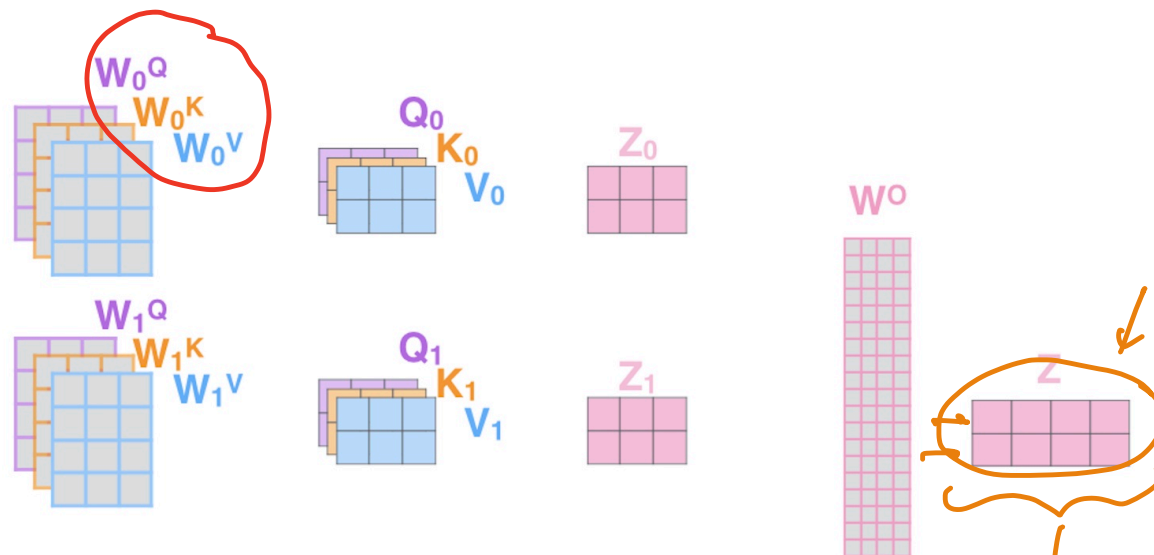* Transformer Implementation using Transfer learning.
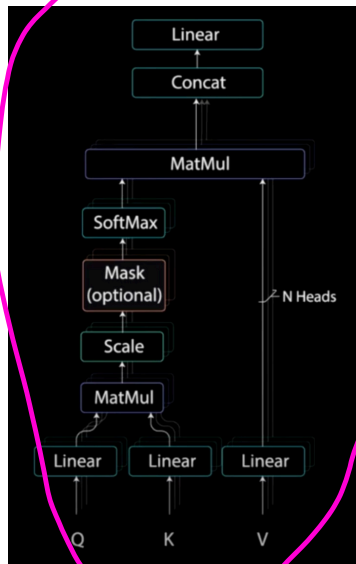
* Using BERT for an NER task.

**X**
Thinking Machines

98

* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one
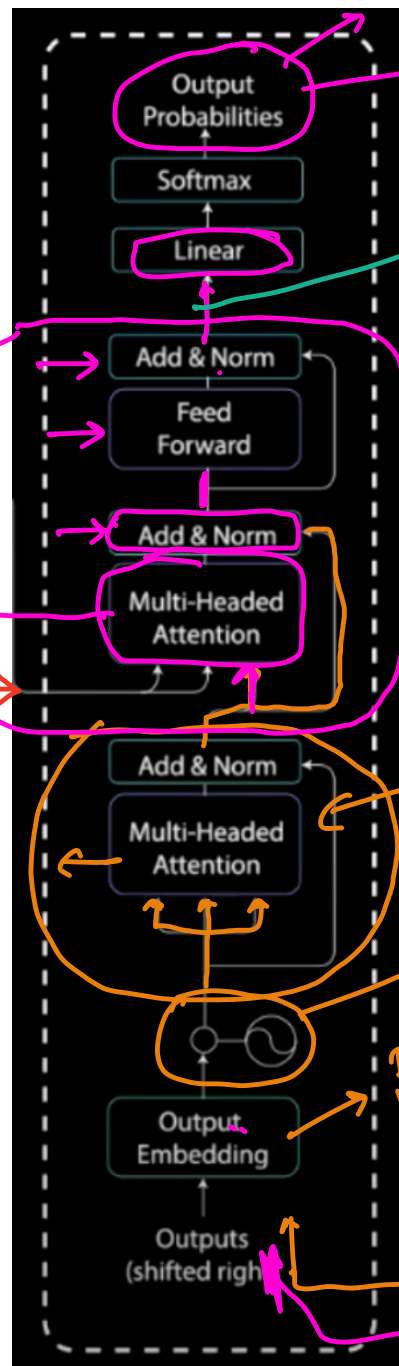
$W_0^Q$
$W_0^K$
$W_0^V$

$W_1^Q$
$W_1^K$
$W_1^V$

$Q_0$
$K_0$
$V_0$

$Q_1$
$K_1$
$V_1$

$Z_0$

$Z_1$

$W^O$

$Z$

→ o/p of multiheaded attention.

Generates 1 word at a time.

New value vector.

⟨start⟩ _0_   This _1_   is _2_   a _3_

cat _4_   ⟨end⟩ _5_

Decoder block. × N.

Encoder Outputs:

✓ Key Matrix
✓ Query Matrix.

Value vectors.

This changes for every new word which is generated.

Positional encoding.

From Glove

$\langle \text{start} \rangle$ $\left\{ \begin{array}{c} \\ \\ \end{array} \right.$

$\left( \dfrac{1}{4} \right)$

$e^{\boxed{-\infty}} = 0 \,?$

$4_{-\infty}$

|   |   |   |   |
|---|---|---|---|
|   |   |   |   |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |

$\left. \begin{array}{c} \\ \\ \end{array} \right\} \dfrac{1}{4}$

↑

**Issue:** Softmax will artificially put $\frac{1}{4}$ as the attention weight when we use padded inputs.

|  | <start> | | | |
|---|---|---|---|---|
| <start> | 0.7 | 0.1 | 0.1 | 0.1 |
| Amo | 0.1 | 0.2 | 0.6 | 0.1 |
| el | 0.1 | 0.6 | 0.2 | 0.1 |
| aprendijaze | 0.1 | 0.2 | 02 | 0.5 |

$$\left( \frac{Q^\top K}{\sqrt{d}} \right)$$

+

| 0 | -inf | -inf | -inf |
|---|---|---|---|
| 0 | 0 | -inf | -inf |
| 0 | 0 | 0 | -inf |
| 0 | 0 | 0 | 0 |

**Mask**

| 0.7 | -inf | -inf | -inf |
|---|---|---|---|
| 0.1 | 0.2 | -inf | -inf |
| 0.1 | 0.6 | 0.2 | -inf |
| 0.1 | 0.2 | 02 | 0.5 |

# Transformers

**BERT.**

## Encoder
→ Study a corpus and store context info.

$\underbrace{\qquad\qquad}$

**Multiple Encoders can be trained for different sources of info.**

## Decoder.

GPT

→ Use context info to answer question.
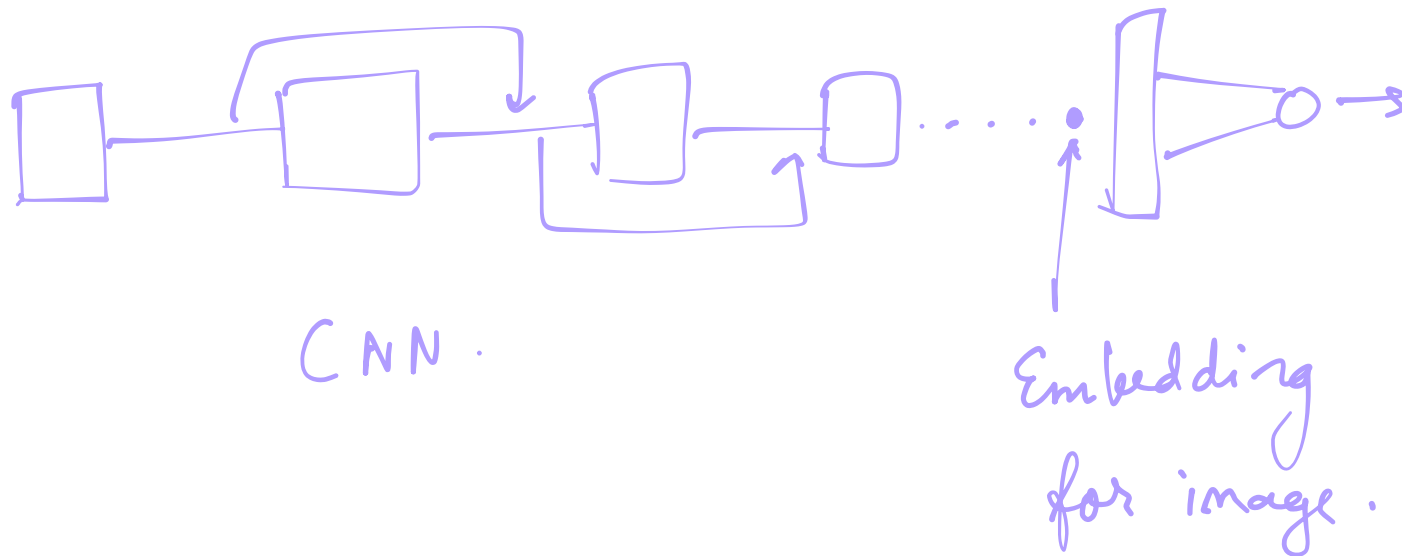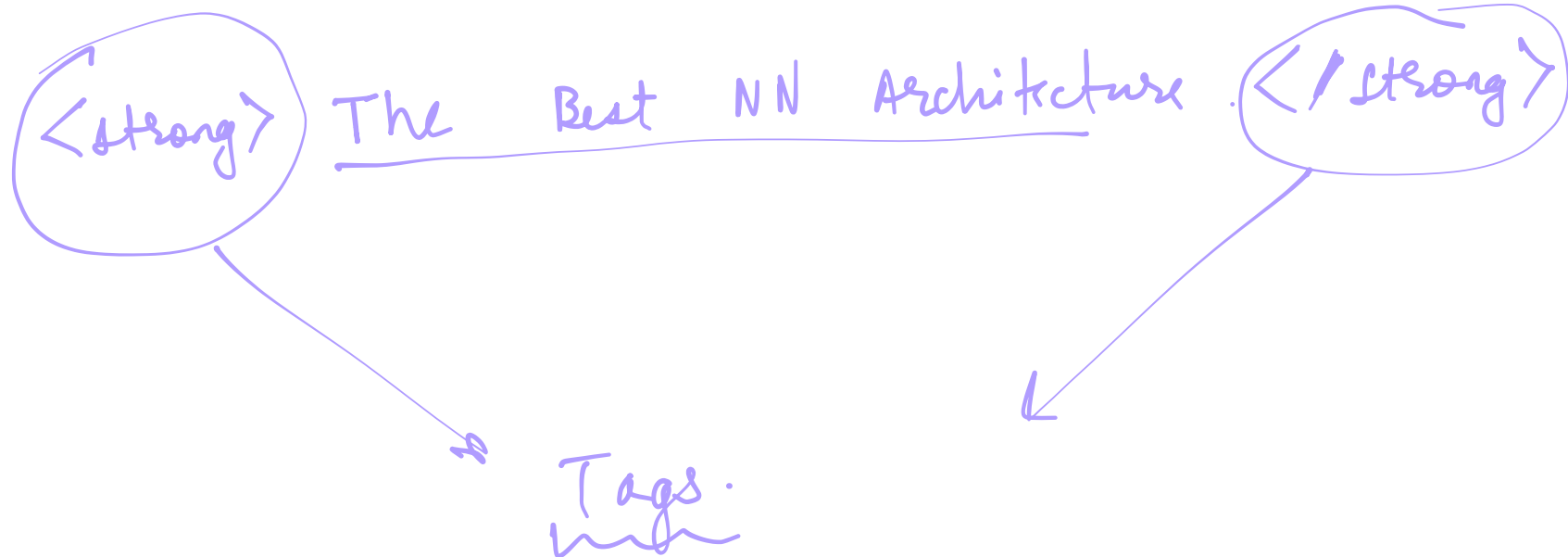
$\underbrace{\qquad\qquad}$ → QnA

Decoder 1 — Task 1 ⌐

Decoder 2 — Task-2 → NER

Decoder 3 — Task 3 → MT.

# Recommender Systems $\longrightarrow$ Embeddings using Matrix factorization.



CNN.

Embedding for image.

**&lt;strong&gt;** The Best NN Architecture . **&lt;/strong&gt;**

Tags.

## Closing Remarks:

* **Heuristic methods.** : Human intervention for useful info from text.

* **Probabilistic methods;** CRF (Conditional Random fields)

   Naïve Bayes.

* **Neural network approaches:**

Conv → ReLU → BN → Conv ×N.

Q, K

Attention.