# Introduction to Deep Learning

## Lecture 19
## Transformers

**Liangze Li**            **Kateryna Shapovalenko**
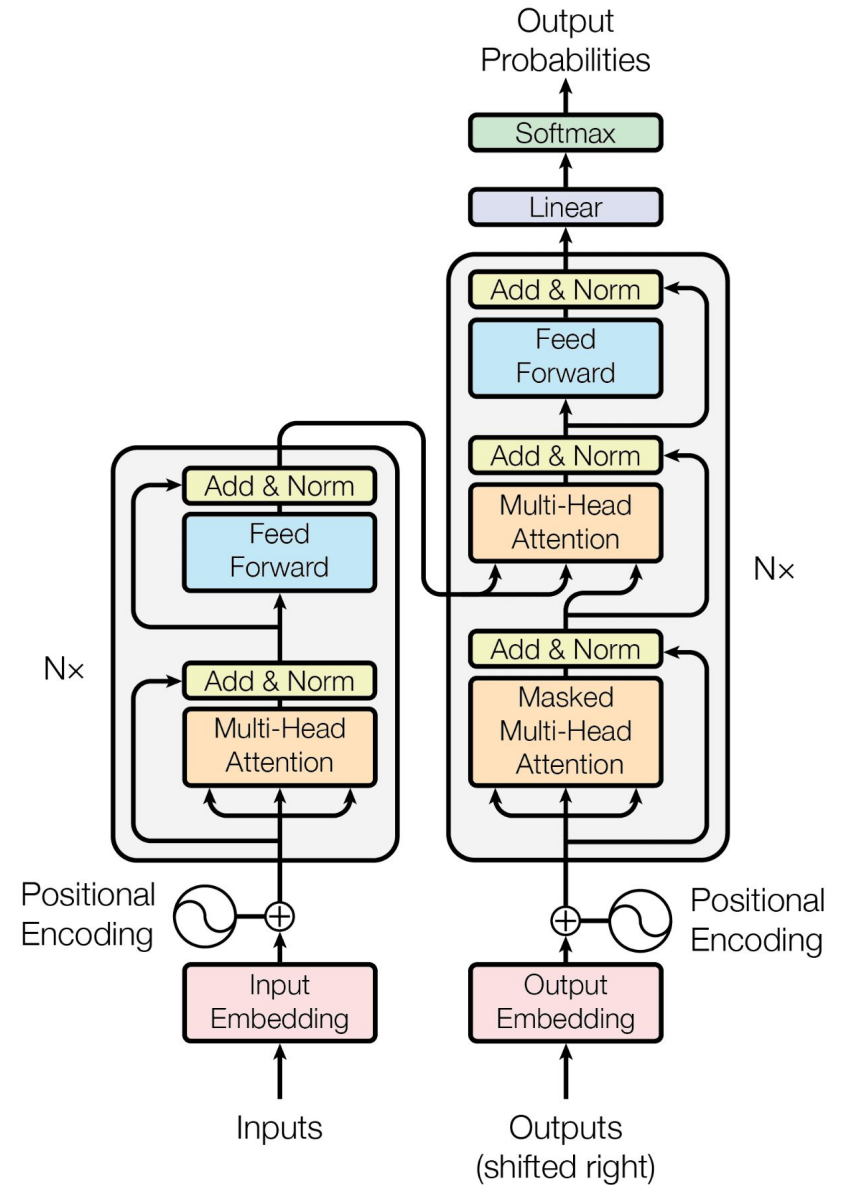
**11-785, Spring 2024**

# Attendance poll @1585

# Table of contents

1. **The Transformer Architecture**

2. **Pre-training and Fine-tuning**

3. **Transformer Applications**

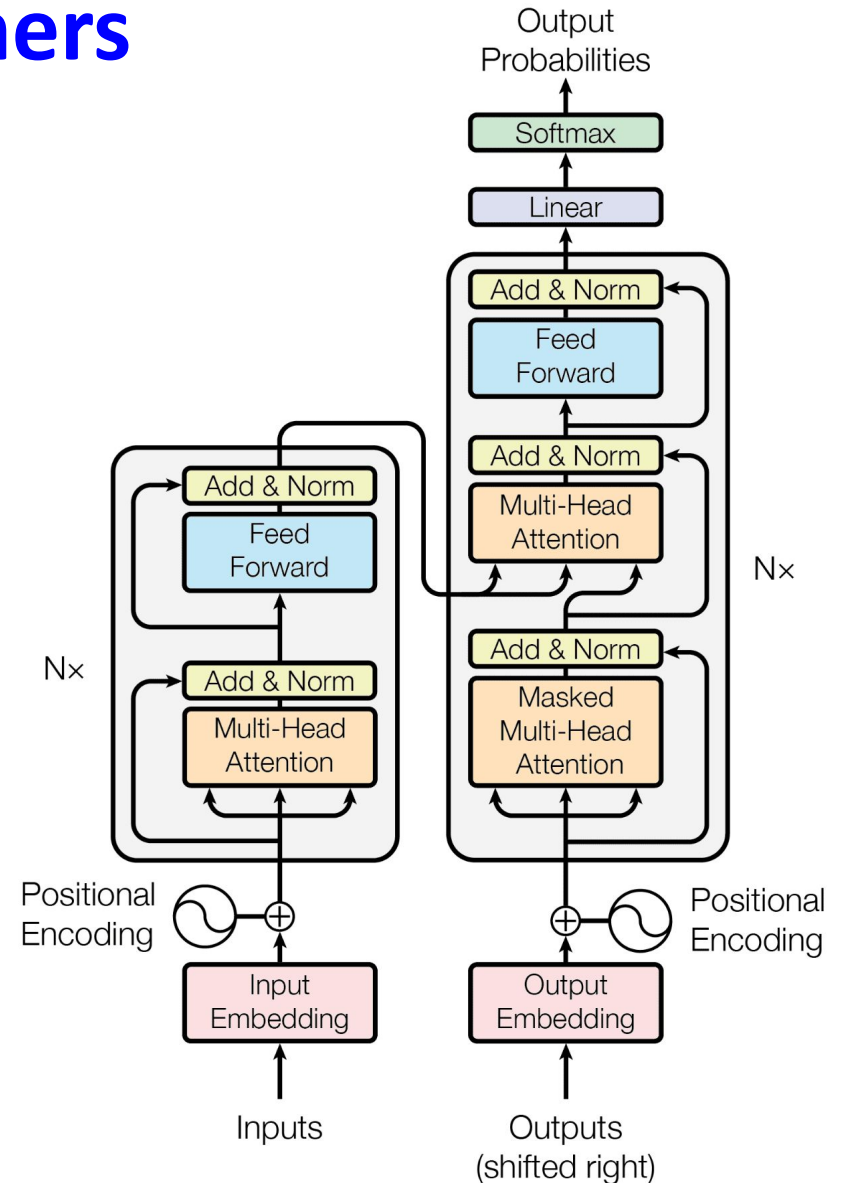4. **Case study - Large Language Models**

# Part 1

# Transformer Architecture

# Transformers

- Tokenization
- Input Embeddings
- Position Encodings
- Query, Key, & Value
- Attention
- Self Attention
- Multi-Head Attention
- Feed Forward
- Add & Norm
- Encoders

- Masked Attention
- Encoder Decoder Attention
- Linear
- Softmax
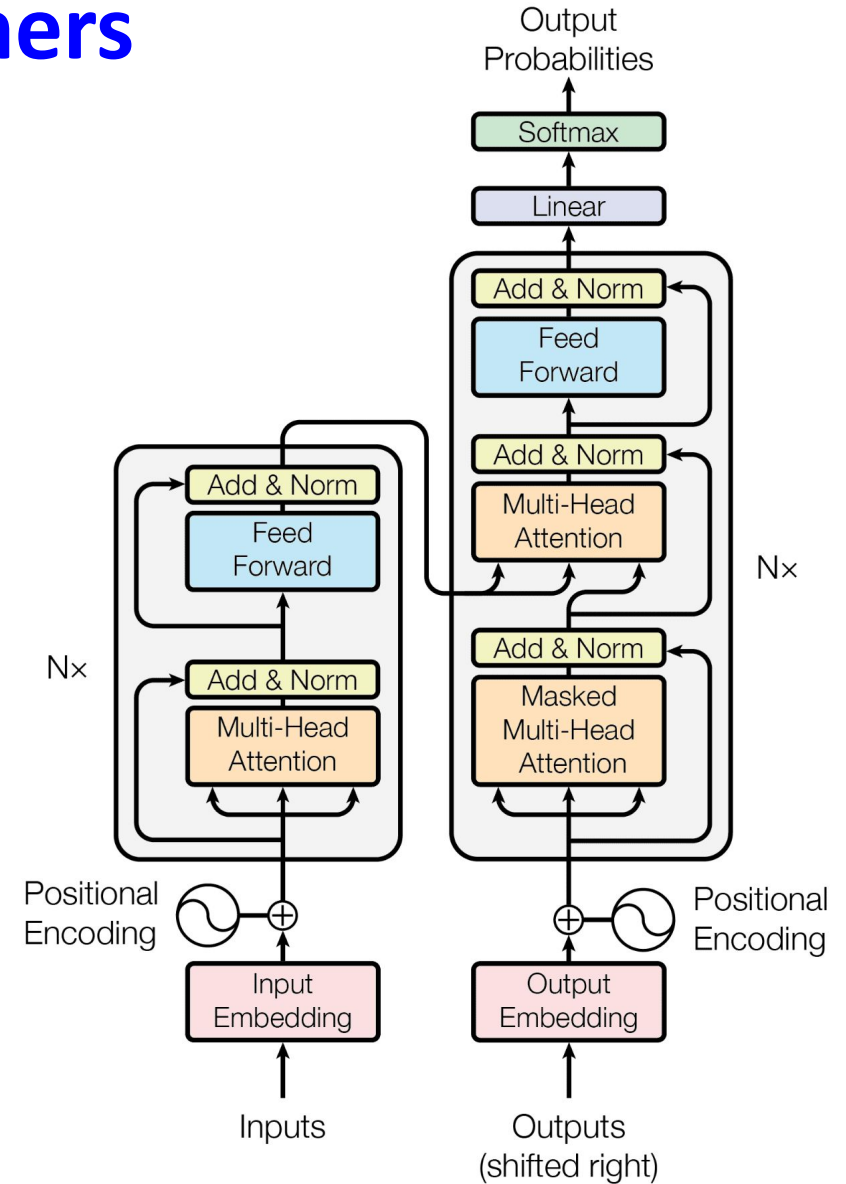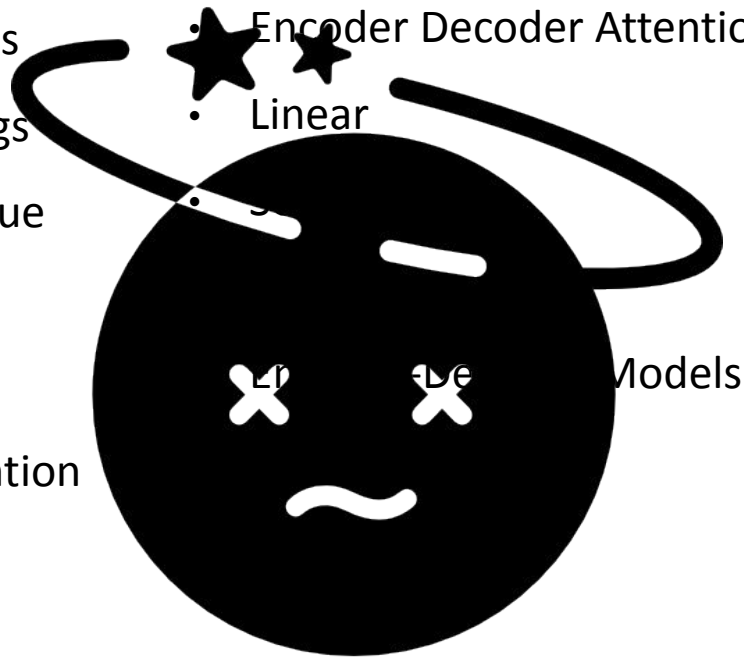- Decoders
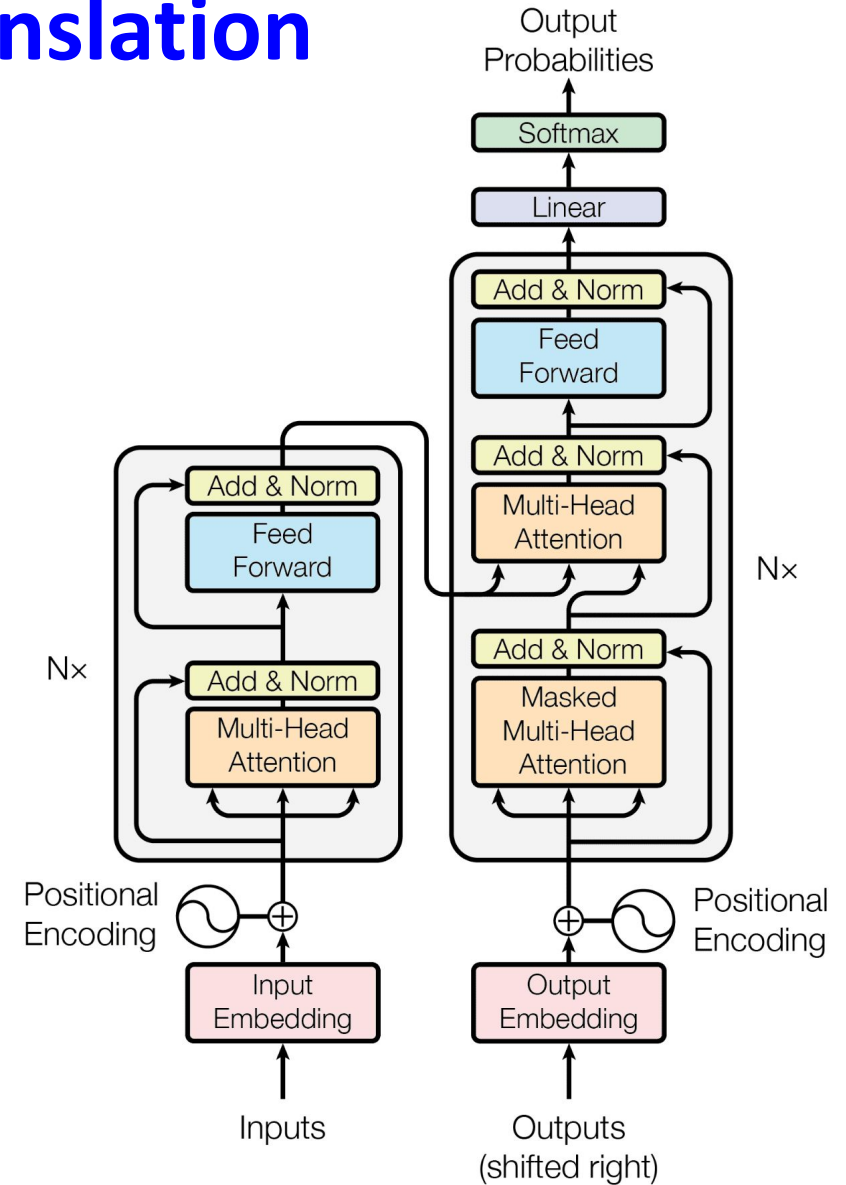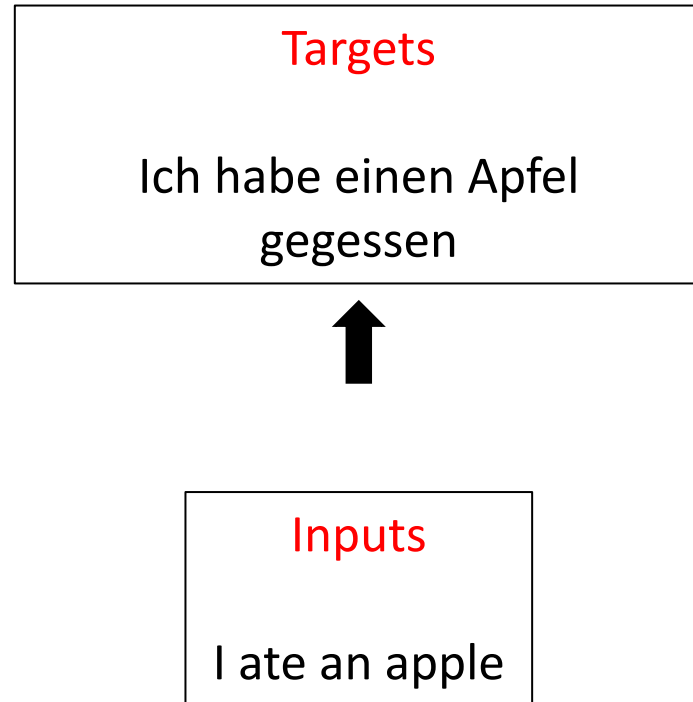- Encoder-Decoder Models

# Transformers

- Tokenization
- Input Embeddings
- Position Encodings
- Query, Key, & Value
- Attention
- Self Attention
- Multi-Head Attention
- Feed Forward
- Add & Norm
- Encoders

- Masked Attention
- Encoder Decoder Attention
- Linear

Models

# **Machine Translation**



Targets

Ich habe einen Apfel gegessen

Inputs

I ate an apple

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Nx

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Nx

Add & Norm

Masked Multi-Head Attention

Positional Encoding

Input Embedding

Positional Encoding

Output Embedding

Inputs

Outputs (shifted right)

# Inputs

## Processing Inputs

| Inputs |
|---|
| I ate an apple |



Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

N×

Add & Norm

Feed
Forward

N×

Add & Norm

Multi-Head
Attention

Add & Norm

Masked
Multi-Head
Attention

Positional
Encoding

Positional
Encoding

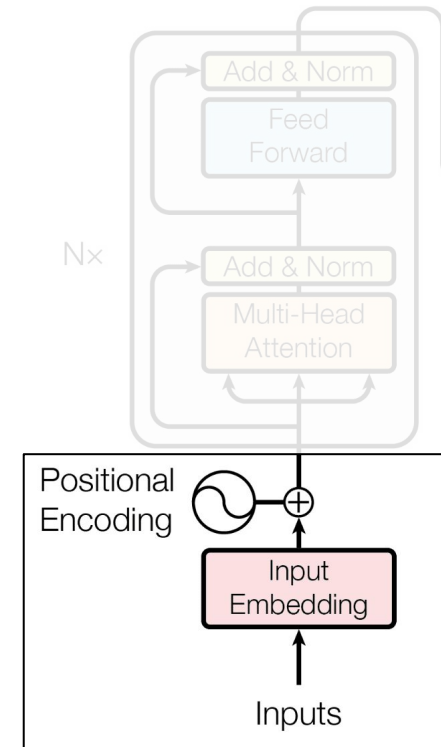Input
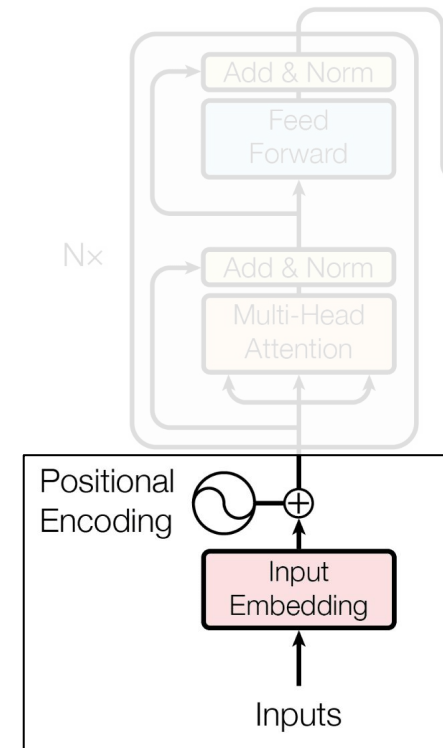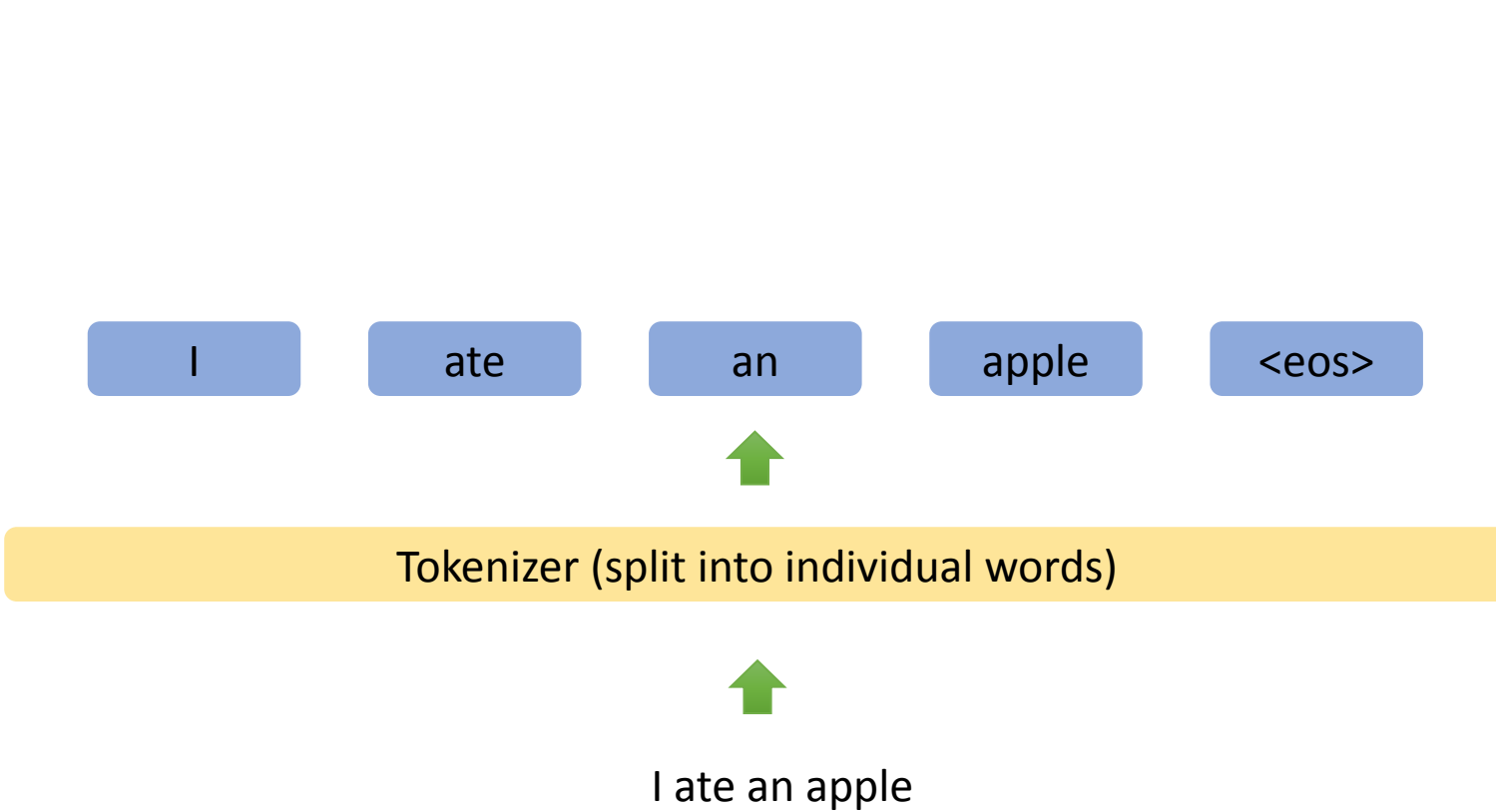Embedding

Output
Embedding

Inputs

Outputs
(shifted right)

# Tokenization

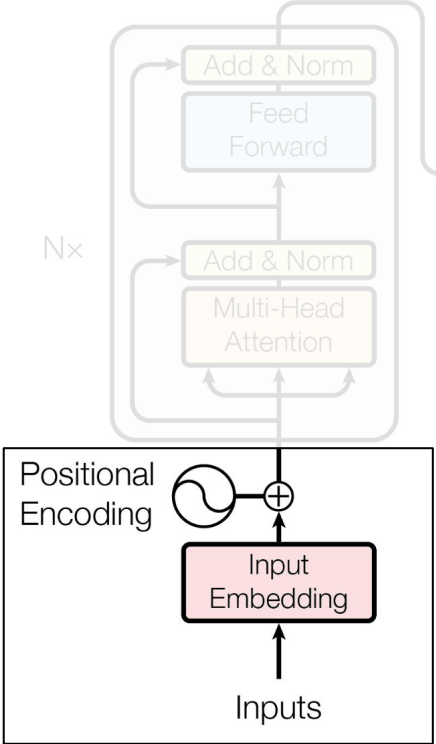Tokenizer (split into individual words)

I ate an apple

Positional Encoding

Input Embedding

Inputs

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Nx

# Tokenization

| I | ate | an | apple | <eos> |

Tokenizer (split into individual words)

I ate an apple

# Input Embeddings

Embedding Layer

| I | ate | an | apple | <eos> |

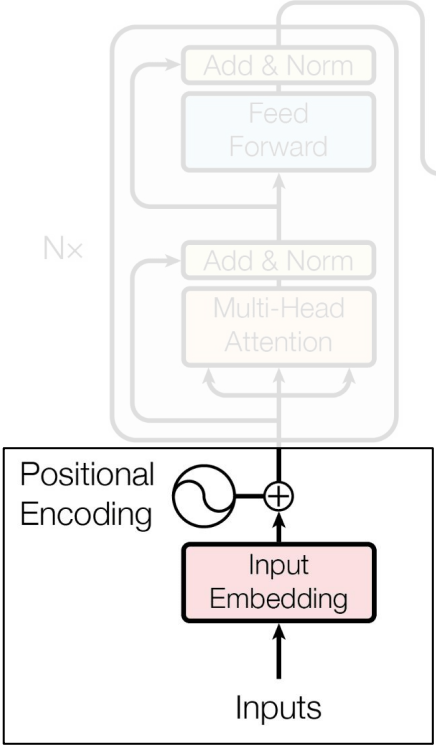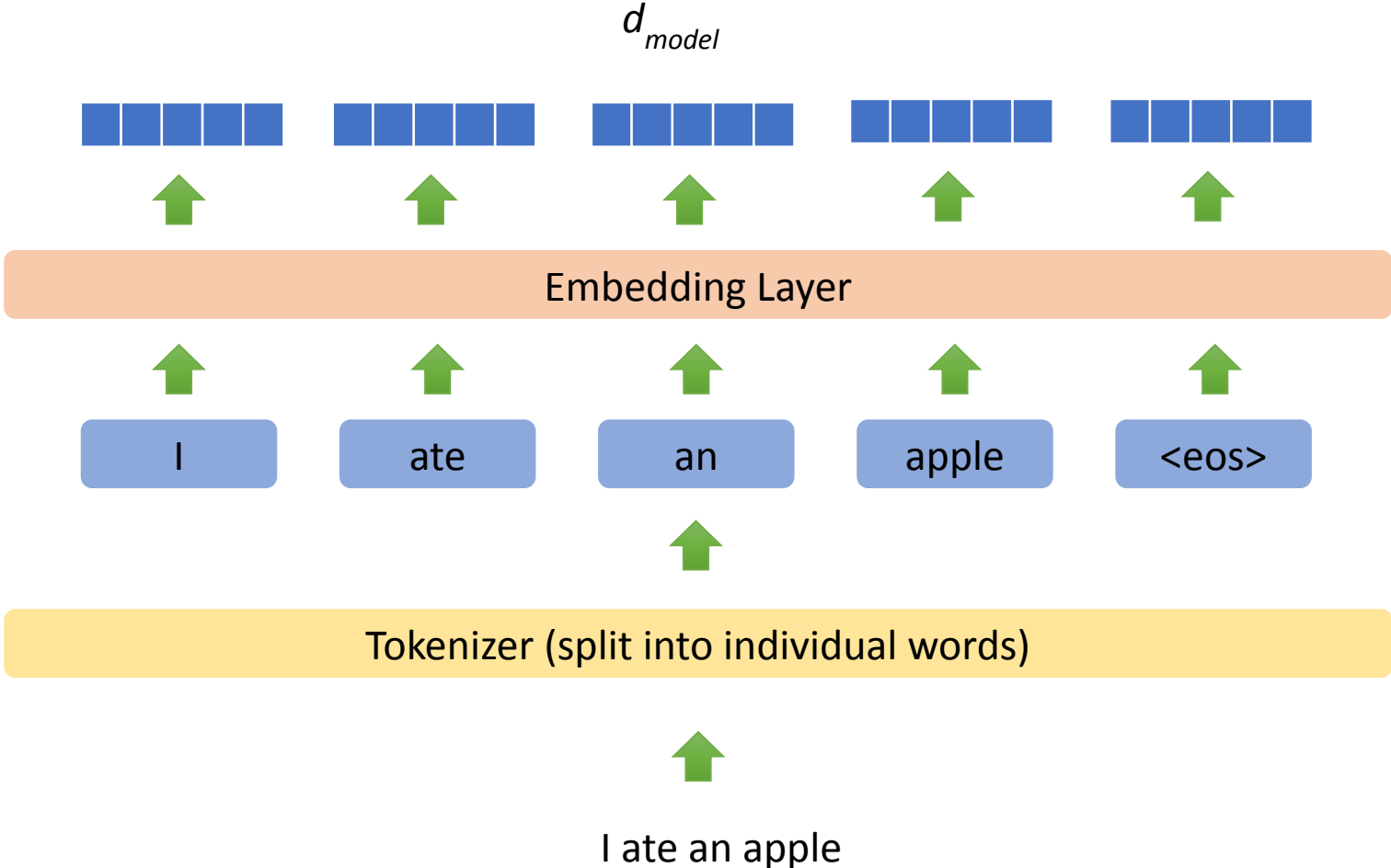Tokenizer (split into individual words)

I ate an apple

**Generate Input Embeddings**

Add & Norm

Feed Forward

Nx

Add & Norm

Multi-Head Attention
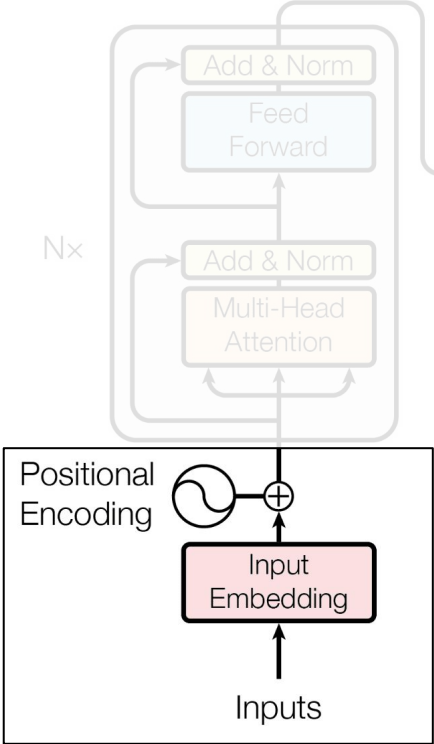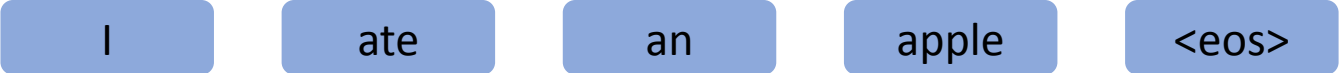
Positional Encoding

Input Embedding

Inputs

# Input Embeddings



Generate Input Embeddings

# Position Encodings

| I | ate | an | apple | <eos> |



Positional Encoding

Input Embedding

Inputs

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

# Position Encodings

| I | ate | an | apple | <eos> |

❌

| apple | ate | an | I | <eos> |

Positional Encoding

Input Embedding

Inputs

Add & Norm

Feed Forward

Nx

Add & Norm

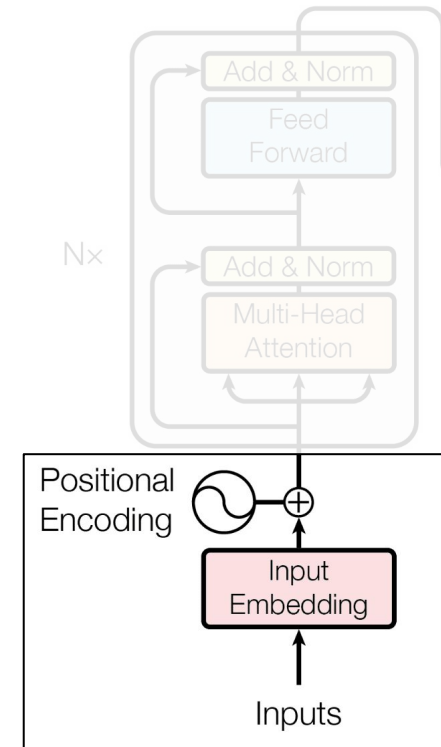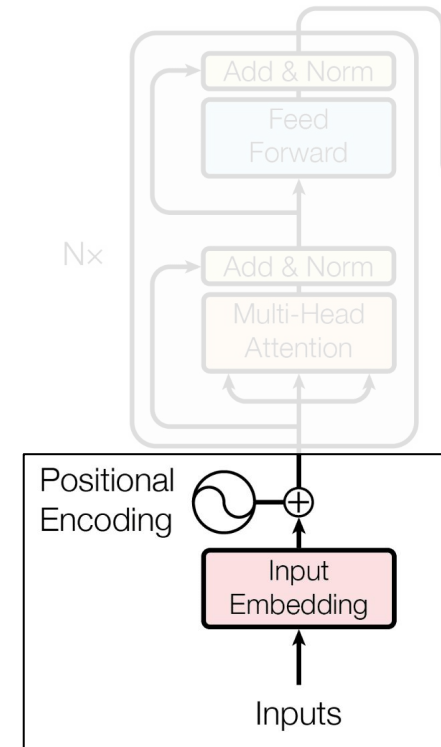Multi-Head Attention

# Position Encodings

**Requirements for Positional Encodings???**

# Position Encodings

**Requirements for Positional Encodings**

- Some representation of time? (like **seq2seq**?)

- Should be unique for each position – not cyclic

# Position Encodings

**Requirements for Positional Encodings**

- Some representation of time? (like **seq2seq**?)

- Should be unique for each position – not cyclic

Possible Candidates :

$$P_{t+1} = P_t + \Delta c$$

$$P_{t+1} = e^{P_t \Delta c}$$

$$P_{t+1} = P_t^{\cdot \ t\Delta c}$$

# Position Encodings

**Requirements for Positional Encodings**

- Some representation of time? (like **seq2seq**?)

- Should be unique for each position – not cyclic

Possible Candidates:

$$P_{t+1} = P_t + \Delta c$$

$$P_{t+1} = e^{P_{t_\Delta} c}$$

$$P_{t+1} = P_t \cdot {}^{t\Delta c}$$



f(x)=50x

f(x)=x³

f(x)=2ˣ



Positional Encoding

Input Embedding

Inputs

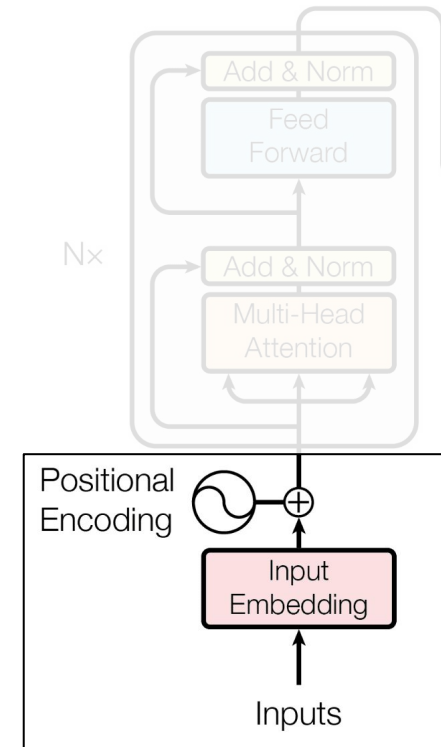# Position Encodings
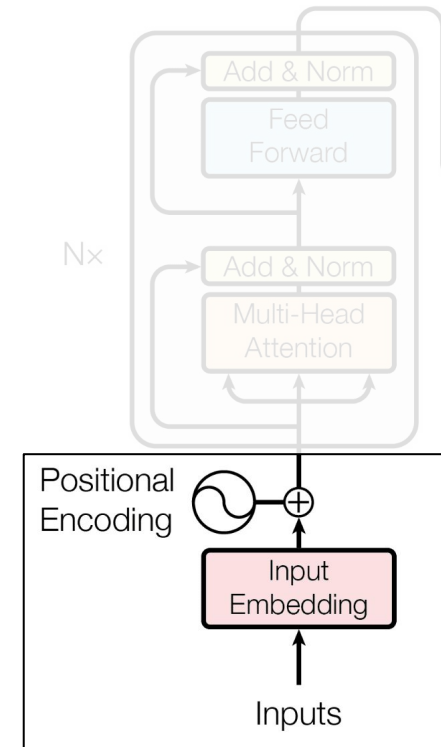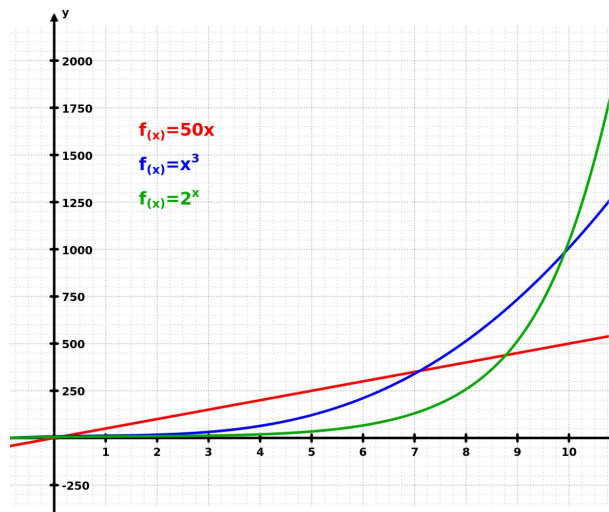
**Requirements for Positional Encodings**

- Some representation of time? (like **seq2seq**?)

- Should be unique for each position – not cyclic

Possible Candidates :

$$P_{t+1} = P_t + \Delta c$$

$$P_{t+1} = P_t \cdot c$$

$$P_{t+1} = P_t \cdot {}^{t\Delta c}$$

$f_{(x)}=50x$

$f_{(x)}=x^3$

$f_{(x)}=2^x$

# Position Encodings

**Requirements for Positional Encodings**

- Some representation of time? (like **seq2seq**?)

- Should be unique for each position – not cyclic

- **Bounded**

**Possible Candidates**

$$P(t + t') = M^{t'} \times P(t)$$

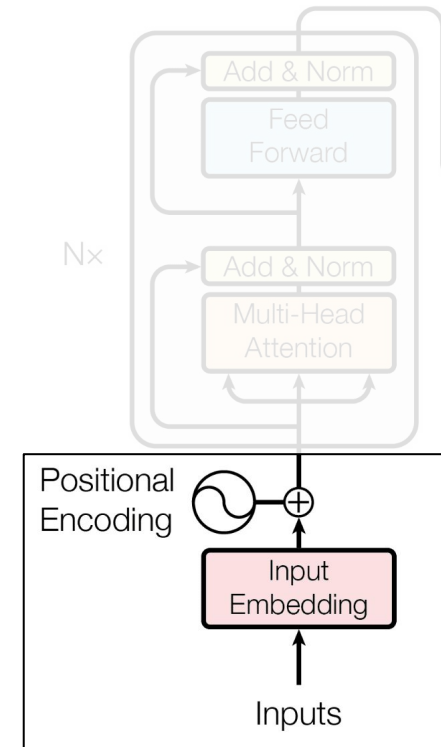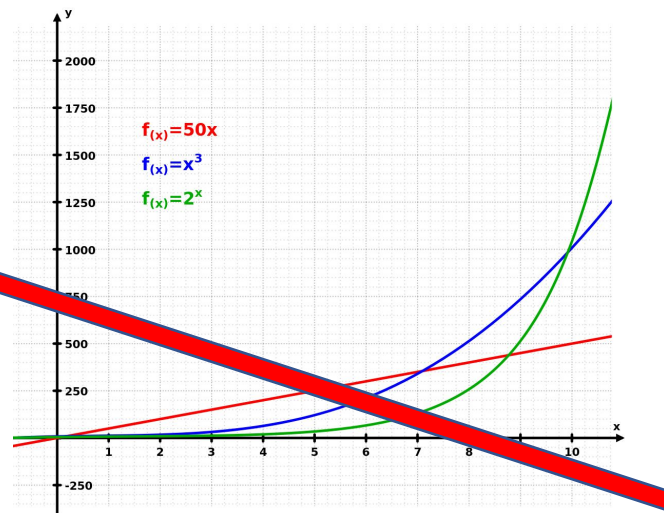# Position Encodings

**Requirements for Positional Encodings**

- Some representation of time? (like **seq2seq**?)

- Should be unique for each position – not cyclic

- **Bounded**

**Possible Candidates**

$$P(t + t') = \mathbf{M}^{t'} \times P(t)$$

**M?**

1. **Should be a unitary matrix**

2. **Magnitudes of eigen value should be 1 -> norm preserving**

3. **The matrix can be learnt**

4. **Produces unique rotated embeddings each time**

# Rotary Position Embedding

## RoFormer: Enhanced Transformer with Rotary Position Embedding

$$f_{\{q,k\}}(\boldsymbol{x}_m, m) = \begin{pmatrix} \cos m\theta & -\sin m\theta \\ \sin m\theta & \cos m\theta \end{pmatrix} \begin{pmatrix} W_{\{q,k\}}^{(11)} & W_{\{q,k\}}^{(12)} \\ W_{\{q,k\}}^{(21)} & W_{\{q,k\}}^{(22)} \end{pmatrix} \begin{pmatrix} x_m^{(1)} \\ x_m^{(2)} \end{pmatrix}$$

Table 2: Comparing RoFormer and BERT by fine tuning on downstream GLEU tasks.

| Model | MRPC | SST-2 | QNLI | STS-B | QQP | MNLI(m/mm) |
|-------|------|-------|------|-------|------|------------|
| BERTDevlin et al. [2019] | 88.9 | 93.5 | 90.5 | 85.8 | 71.2 | 84.6/83.4 |
| RoFormer | **89.5** | 90.7 | 88.0 | **87.0** | **86.4** | 80.2/79.8 |



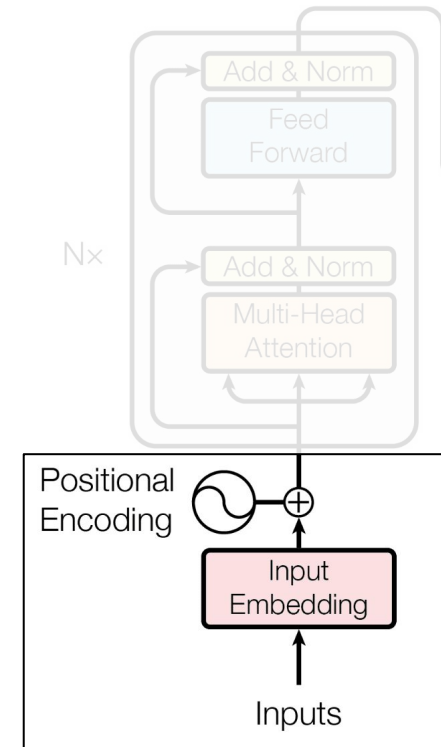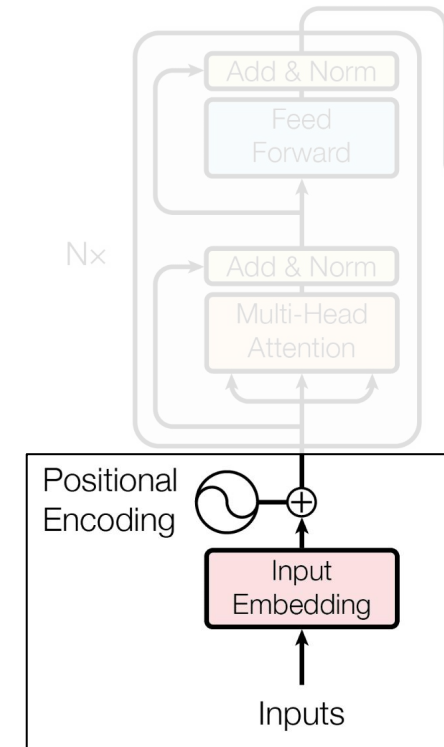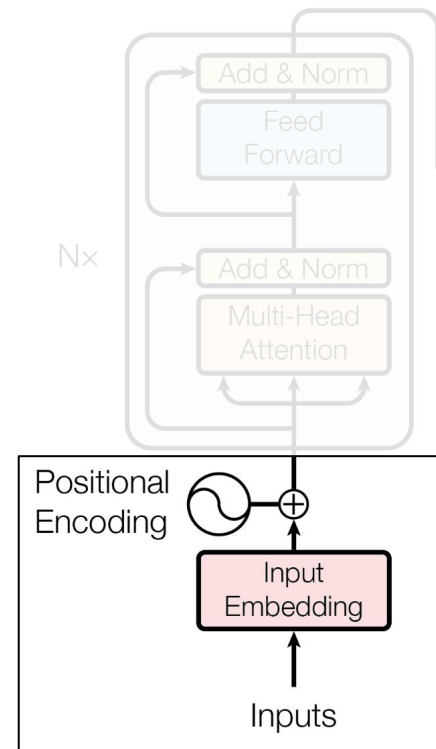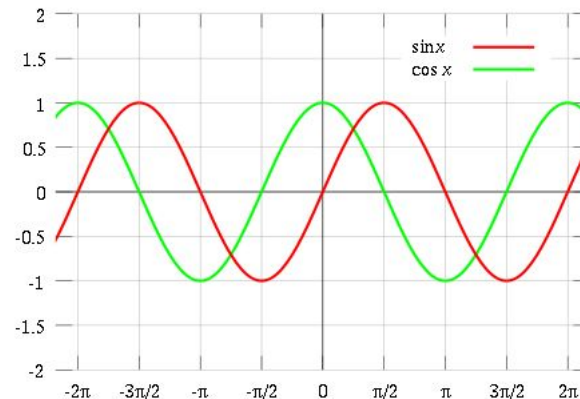REF: Rotary Position Embeddings 🔗

# Position Encoding

**Requirements for Position Encodings**

- Some representation of time? (like **seq2seq**?)

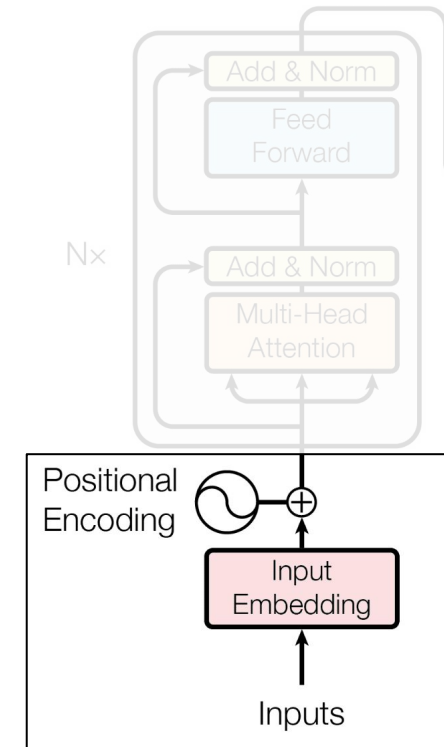- Should be unique for each position

- Bounded

**Actual Candidates**

    *sine(g(t))*

    *cosine(g(t))*

**Requirements for g(t)**

- Must have same dimensions as input embeddings

- Must produce overall unique encodings

# Position Encoding

For each position, an embedded input is moved the same distance but at a different angle. **Inputs that are close to each other in the sequence have similar perturbations, but inputs that are far apart are perturbed in different directions**.

$$PE_{(pos,2i)} = sin(pos/10000^{2i/d_{model}})$$

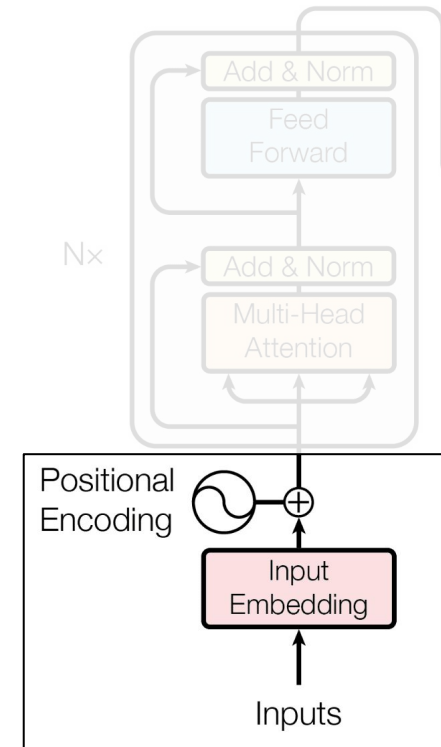$$PE_{(pos,2i+1)} = cos(pos/10000^{2i/d_{model}})$$

**pos** -> idx of the token in input sentence

**i** -> $i^{th}$ dimension out of d

**d model** -> embedding dimension of each token

Different calculations for odd and even embedding indices



24

# Position Encoding
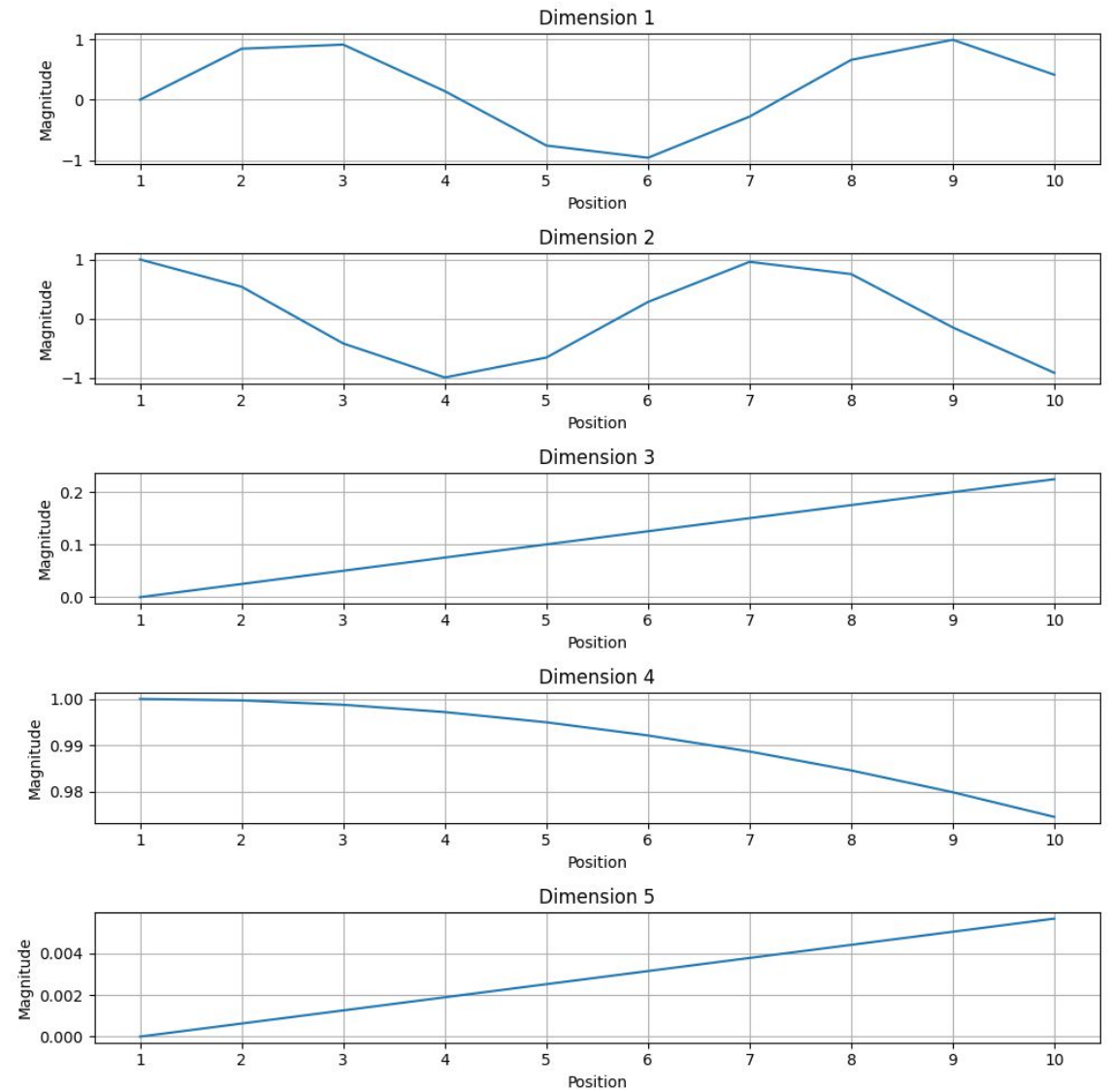


$$PE_{(pos,2i)} = sin(pos/10000^{2i/d_{model}})$$

$$PE_{(pos,2i+1)} = cos(pos/10000^{2i/d_{model}})$$

```
Positional Encoding:
          0      1       2       3      4
Dim 1 0.000  0.841   0.909   0.141 -0.757
Dim 2 1.000  0.540  -0.416  -0.990 -0.654
Dim 3 0.000  0.025   0.050   0.075  0.100
Dim 4 1.000  1.000   0.999   0.997  0.995
Dim 5 0.000  0.001   0.001   0.002  0.003
```

# Position Encoding



Final Input Embeddings

Position Encodings

Input Embeddings

Embedding Layer

| I | ate | an | apple | <eos> | Tokens

Tokenizer

I ate an apple

Input

Positional Encoding

Input Embedding

Inputs

# Transformers

✔ **Tokenization**

✔ **Input Embeddings**

✔ **Position Encodings**

- Query, Key, & Value

- Attention

- Self Attention

- Multi-Head Attention

- Feed Forward

- Add & Norm

- Encoders

- Masked Attention

- Encoder Decoder Attention

- Linear

- Softmax

- Decoders

- Encoder-Decoder Models

# Encoder

WHERE IS THE CONTEXT ?



| | | | | |
|---|---|---|---|---|
| I | ate | an | apple | <eos> |

# Encoder

**BLACK BOX OF SORTS**



Add & Norm

Feed Forward

N×

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

| I | ate | an | apple | <eos> |

# Encoder

BLACK BOX
OF SORTS

**LEARN TO ADD CONTEXT**



I    ate    an    apple    <eos>

Positional Encoding

Input Embedding

Inputs

Nx

Add & Norm
Feed Forward
Add & Norm
Multi-Head Attention

# Encoder

CONTEXTUALLY RICH EMBEDDINGS

BLACK BOX
OF SORTS

LEARN TO
ADD
CONTEXT

I    ate    an    apple    <eos>

Add & Norm

Feed
Forward

N×    Add & Norm

Multi-Head
Attention

Positional
Encoding

Input
Embedding

Inputs

# Encoder

$$\alpha_{[ij]} \ ?$$

**CONTEXTUALLY RICH EMBEDDINGS**

**BLACK BOX OF SORTS**

**LEARN TO ADD CONTEXT**

I  ate  an  apple  <eos>

Add & Norm

Feed Forward

N×

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

# Attention

$$\alpha_{[ij]} \ ?$$

From lecture 18:

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

# Attention

$$\alpha_{[ij]} \; ?$$

From lecture 18:

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

- Query

- Key

- Value



34

# Query, Key & Value

**Database**

**{Key, Value store}**
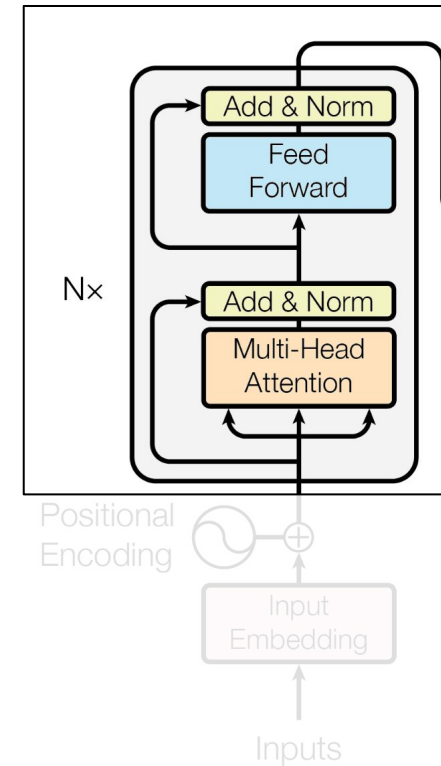
{"order_100": {"items":"a1", "delivery_date":"a2", ...}},
{"order_101": {"items":"b1", "delivery_date":"b2", ...}},
{"order_102": {"items":"c1", "delivery_date":"c2", ...}},
{"order_103": {"items":"d1", "delivery_date":"d2", ...}},
{"order_104": {"items":"e1", "delivery_date":"e2", ...}},
{"order_105": {"items":"f1", "delivery_date":"f2", ...}},
{"order_106": {"items":"g1", "delivery_date":"g2", ...}},
{"order_107": {"items":"h1", "delivery_date":"h2", ...}},
{"order_108": {"items":"i1", "delivery_date":"i2", ...}},
{"order_109": {"items":"j1", "delivery_date":"j2", ...}},
{"order_110": {"items":"k1", "delivery_date":"k2", ...}}

# Query, Key & Value

**Database**

**{Key, Value store}**

{Query: "Order details of order_104"}

OR

{Query: "Order details of order_106"}

{"order_100": {"items":"a1", "delivery_date":"a2", ....}},
{"order_101": {"items":"b1", "delivery_date":"b2", ...}},
{"order_102": {"items":"c1", "delivery_date":"c2", ...}},
{"order_103": {"items":"d1", "delivery_date":"d2", ...}},
{"order_104": {"items":"e1", "delivery_date":"e2", ...}},
{"order_105": {"items":"f1", "delivery_date":"f2", ...}},
{"order_106": {"items":"g1", "delivery_date":"g2", ...}},
{"order_107": {"items":"h1", "delivery_date":"h2", ...}},
{"order_108": {"items":"i1", "delivery_date":"i2", ...}},
{"order_109": {"items":"j1", "delivery_date":"j2", ...}},
{"order_110": {"items":"k1", "delivery_date":"k2", ....}}

# Query, Key & Value

{Key, Value store}

{Query: "Order details of order_104"}

OR

{Query: "Order details of order_106"}

```
{"order_100": {"items":"a1", "delivery_date":"a2", ....}},
{"order_101": {"items":"b1", "delivery_date":"b2", ...}},
{"order_102": {"items":"c1", "delivery_date":"c2", ...}},
{"order_103": {"items":"d1", "delivery_date":"d2", ...}},
{"order_104": {"items":"e1", "delivery_date":"e2", ...}},
{"order_105": {"items":"f1", "delivery_date":"f2", ...}},
{"order_106": {"items":"g1", "delivery_date":"g2", ...}},
{"order_107": {"items":"h1", "delivery_date":"h2", ...}},
{"order_108": {"items":"i1", "delivery_date":"i2", ...}},
{"order_109": {"items":"j1", "delivery_date":"j2", ...}},
{"order_110": {"items":"k1", "delivery_date":"k2", ...}}
```

# Query, Key & Value

{Key, Value store}

{Query: "Order details of order_104"}

OR

{Query: "Order details of order_106"}

```
{"order_100": {"items":"a1", "delivery_date":"a2", .....}},
{"order_101": {"items":"b1", "delivery_date":"b2", ...}},
{"order_102": {"items":"c1", "delivery_date":"c2", ...}},
{"order_103": {"items":"d1", "delivery_date":"d2", ...}},
{"order_104": {"items":"e1", "delivery_date":"e2", ...}},
{"order_105": {"items":"f1", "delivery_date":"f2", ...}},
{"order_106": {"items":"g1", "delivery_date":"g2", ...}},
{"order_107": {"items":"h1", "delivery_date":"h2", ...}},
{"order_108": {"items":"i1", "delivery_date":"i2", ...}},
{"order_109": {"items":"j1", "delivery_date":"j2", ...}},
{"order_110": {"items":"k1", "delivery_date":"k2", ...}}
```

# Query, Key & Value

{Key, Value store}

{Query: "Order details of order_104"}

OR

{Query: "Order details of order_106"}

{"order_100": {"items":"a1", "delivery_date":"a2", ....}},
{"order_101": {"items":"b1", "delivery_date":"b2", ...}},
{"order_102": {"items":"c1", "delivery_date":"c2", ...}},
{"order_103": {"items":"d1", "delivery_date":"d2", ...}},
{"order_104": {"items":"e1", "delivery_date":"e2", ...}},
{"order_105": {"items":"f1", "delivery_date":"f2", ...}},
{"order_106": {"items":"g1", "delivery_date":"g2", ...}},
{"order_107": {"items":"h1", "delivery_date":"h2", ...}},
{"order_108": {"items":"i1", "delivery_date":"i2", ...}},
{"order_109": {"items":"j1", "delivery_date":"j2", ...}},
{"order_110": {"items":"k1", "delivery_date":"k2", ...}}

# Query, Key & Value

Done at the same time !!

{Key, Value store}

{Query: "Order details of order_104"}

OR

{Query: "Order details of order_106"}

{"order_100": {"items":"a1", "delivery_date":"a2", ....}},
{"order_101": {"items":"b1", "delivery_date":"b2", ...}},
{"order_102": {"items":"c1", "delivery_date":"c2", ...}},
{"order_103": {"items":"d1", "delivery_date":"d2", ...}},
{"order_104": {"items":"e1", "delivery_date":"e2", ...}},
{"order_105": {"items":"f1", "delivery_date":"f2", ...}},
{"order_106": {"items":"g1", "delivery_date":"g2", ...}},
{"order_107": {"items":"h1", "delivery_date":"h2", ...}},
{"order_108": {"items":"i1", "delivery_date":"i2", ...}},
{"order_109": {"items":"j1", "delivery_date":"j2", ...}},
{"order_110": {"items":"k1", "delivery_date":"k2", ...}}

# Query, Key & Value

{Query: "Order details of order_104"}

OR

{Query: "Order details of order_106"}

{"order_100": {"items":"a1", "delivery_date":"a2", ...}},
{"order_101": {"items":"b1", "delivery_date":"b2", ...}},
{"order_102": {"items":"c1", "delivery_date":"c2", ...}},
{"order_103": {"items":"d1", "delivery_date":"d2", ...}},
{"order_104": {"items":"e1", "delivery_date":"e2", ...}},
{"order_105": {"items":"f1", "delivery_date":"f2", ...}},
{"order_106": {"items":"g1", "delivery_date":"g2", ...}},
{"order_107": {"items":"h1", "delivery_date":"h2", ...}},
{"order_108": {"items":"i1", "delivery_date":"i2", ...}},
{"order_109": {"items":"j1", "delivery_date":"j2", ...}},
{"order_110": {"items":"k1", "delivery_date":"k2", ...}}

## Query
1. Search for info

## Key
1. Interacts directly with Queries
2. Distinguishes one object from another
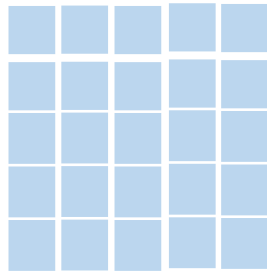3. Identify which object is the most relevant and by how much

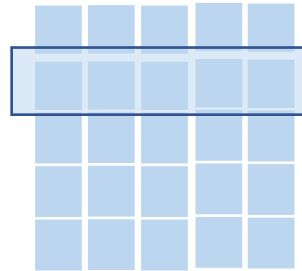## Value
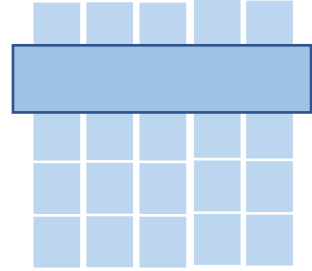1. Actual details of the object
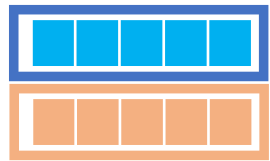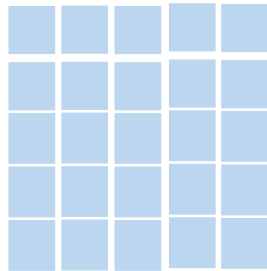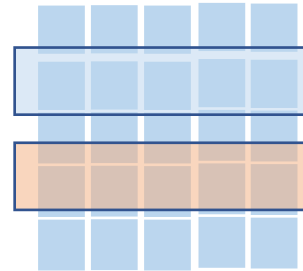2. More fine grained

41

# Attention

Query

Key Value
Store

Key

Value

# Attention



Query

Key Value
Store

Key

Value

# Attention



Done at the same time !!

Query

Key Value Store

Key

Value

44

# Attention

**Parallelizable !!!**

| Query | Key Value Store | Key | Value |
|:---:|:---:|:---:|:---:|

$$Q \qquad QK^T \qquad softmax(\frac{QK^T}{\sqrt{d}}) \qquad softmax(\frac{QK^T}{\sqrt{d}})V$$

# Attention

I1

I2

I3

I4

I5

I

ate

an

apple

# Attention

Dimensions across QKV have been dropped for brevity



47

# Attention

Dimensions across QKV have been dropped for brevity



48

# Attention

Dimensions across QKV have been dropped for brevity



I  ate  an  apple  <eos>

49

# Attention

Dimensions across QKV have been dropped for brevity

softmax

$\alpha_{1,1}$

$e_{1,1}$

$Q_1$ $K_1$ $V_1$
$W_Q$ $W_K$ $W_V$
$I_1$

$Q_2$ $K_2$ $V_2$
$W_Q$ $W_K$ $W_V$
$I_2$

$Q_3$ $K_3$ $V_3$
$W_Q$ $W_K$ $W_V$
$I_3$

$Q_4$ $K_4$ $V_4$
$W_Q$ $W_K$ $W_V$
$I_4$

$Q_5$ $K_5$ $V_5$
$W_Q$ $W_K$ $W_V$
$I_5$

I          ate          an          apple          <eos>

Dimensions across QKV have been dropped for brevity

# Attention

# Attention

Dimensions across QKV have been dropped for brevity

# Attention

Dimensions across QKV have been dropped for brevity



53

# Attention

# Attention

Contextually rich embedding

$Z_1$

Dimensions across QKV have been dropped for brevity

softmax

$\alpha_{1,1}$  $\alpha_{1,2}$  $\alpha_{1,3}$  $\alpha_{1,4}$  $\alpha_{1,5}$

$e_{1,1}$  $e_{1,2}$  $e_{1,3}$  $e_{1,4}$  $e_{1,5}$

$Q_1$ $K_1$ $V_1$  $Q_2$ $K_2$ $V_2$  $Q_3$ $K_3$ $V_3$  $Q_4$ $K_4$ $V_4$  $Q_5$ $K_5$ $V_5$

$W_Q$ $W_K$ $W_V$  $W_Q$ $W_K$ $W_V$  $W_Q$ $W_K$ $W_V$  $W_Q$ $W_K$ $W_V$  $W_Q$ $W_K$ $W_V$

$I_1$  $I_2$  $I_3$  $I_4$  $I_5$

I       ate      an      apple      <eos>

55

# Attention



Contextually rich embedding

Dimensions across QKV have been dropped for brevity

Parallelized

56

# Transformers

✔ **Tokenization**

✔ **Input Embeddings**

✔ **Position Encodings**

✔ **Query, Key, & Value**

✔ **Attention**

- Self Attention

- Multi-Head Attention

- Feed Forward

- Add & Norm

- Encoders

- Masked Attention

- Encoder Decoder Attention

- Linear

- Softmax

- Decoders

- Encoder-Decoder Models

# Poll 1 - @1581

**Which of the following are true about attention?**

a. To calculate attention weights for input $l_2$, you would use key $k_2$, and all queries

b. To calculate attention weights for input $l_2$, you would use query $q_2$, and all keys

c. We scale the $QK^T$ product to bring attention weights in the range of [0,1]

d. We scale the $QK^T$ product to allow for numerical stability

# Poll 1 - @1581

**Which of the following are true about attention?**

a.  To calculate attention weights for input $I_2$, you would use key $k_2$, and all queries

b.  **To calculate attention weights for input $I_2$, you would use query $q_2$, and all keys**

c.  We scale the $QK^T$ product to bring attention weights in the range of [0,1]

d.  **We scale the $QK^T$ product to allow for numerical stability**

# Self Attention

From lecture 18:

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

# Self Attention

The animal didn't cross the street because it was too wide

# Self Attention

The animal didn't cross the street because **it** was too wide

?

coreference resolution?

# Self Attention

# Self Attention



The animal didn't cross the street beacuse it was too wide

# Self Attention



The | animal | didn't | cross | the | street | beacuse | it | was | too | wide

# Self Attention



**SELF**

Query Inputs    =    Key Inputs    =    Value Inputs

# Self Attention

$$R^{d_{model} \times d_{model}}$$

$W_Q$      $W_K$      $W_v$

$$R^{T \times d_{model}}$$

Input Embeddings

# Self Attention

$R^{T \times d_{model}}$   $R^{d_{model} \times d_{model}}$



Input Embeddings $\otimes$ W$_Q$ $\longrightarrow$ Q Projections   $R^{T \times d_{model}}$

Input Embeddings $\otimes$ W$_K$ $\longrightarrow$ K Projections   $R^{T \times d_{model}}$

Input Embeddings $\otimes$ W$_V$ $\longrightarrow$ V Projections   $R^{T \times d_{model}}$

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Positional Encoding

Input Embedding

Inputs

68

# Self Attention

softmax $\left[ \dfrac{R^{T \times d_{model}} \underset{Q_{\text{Projection}}}{\square} \otimes R^{d_{model} \times T} \underset{K_{\text{Projection}}}{\square}}{\sqrt{d_{model}}} \right]$ T



$R^{T \times T}$

# Self Attention

$$o\left(T^2 \times d_{model}\right)$$

softmax

$R^{T \times d_{model}}$      $R^{d_{model} \times T}$    T



$Q_{Projection}$        $K_{Projection}$

$$\sqrt{d_{model}}$$

$$R^{T \times T}$$

N×

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Positional
Encoding

Input
Embedding

Inputs

# Self Attention



$$O\left(T^2 \times d_{model}\right)$$

softmax

$R^{T \times d_{model}}$   $R^{d_{model} \times T}$   $R^{T \times d_{model}}$

$Q_{Projection}$   $\otimes$   $K_{Projection}$   $\otimes$   $V_{Projection}$

T

$$\frac{}{\sqrt{d_{model}}}$$

$R^{T \times T}$

Add & Norm

Feed Forward

N×   Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

# Self Attention

$$R^{T \times d_{model}}$$

Attention: Z

# Self Attention



| The | animal | didn't | cross | the | street | because | it | was | too | wide |

Coreference resolution ✓

# Self Attention



The | animal | didn't | cross | the | street | because | it | was | too | wide

Sentence boundaries ?            Coreference resolution ✓            Context ?

Semantic relationships ?            Part of Speech ?            Comparisons ?

# Self Attention

$$R^{d_{model} \times d_{model}}$$



W$_Q$

W$_K$

W$_v$

$$R^{T \times d_{model}}$$

Input Embeddings

# Multi-Head Attention

$$R^{d_{model} \times d_h}$$



H
2
1

$W_{Q1,}$ $W_{Q2, ...}$ $W_{QH,}$

H
2
1

$W_{K1,}$ $W_{K2, ...}$ $W_{KH,}$

H
2
1

$W_{V1,}$ $W_{V2, ...}$ $W_{VH,}$

Input Embeddings

$$R^{T \times d_{model}}$$

$$d_h = \frac{d_{model}}{h}$$



N×

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

# Multi-Head Attention



$R^{T \times d_{model}}$

$R^{d_{model} \times d_h}$

$R^{T \times d_h}$

$R^{T \times d_h}$

$R^{T \times d_h}$

Inputs

Inputs

Inputs

$W_{Qi}$   $Q_i$

$W_{Ki}$   $K_i$

$W_{Vi}$   $V_i$

N×

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

77

# Multi-Head Attention

softmax

$$R^{T \times d_h} \qquad R^{d_h \times T} \qquad\qquad R^{T \times d_h}$$

$$\frac{Q_i \otimes K_i}{\sqrt{d_{model}}} \otimes V_i$$

$$R^{T \times T}$$

for all i ∈ [1, h]

# Multi-Head Attention

$R^{T \times d_h}$     $R^{T \times d_h}$     $R^{T \times d_h}$

$Z_1$     $Z_2$     ...     $Z_h$

CONCAT

Multi Head Attention : Z

$$d_h = \frac{d_{model}}{h}$$

$R^{T \times d_{model}}$

Add & Norm

Feed
Forward

Add & Norm

N×

Multi-Head
Attention

Positional
Encoding

Input
Embedding

Inputs

# Multi-Head
# Attention

The | animal | didn't | cross | the | street | because | it | was | too | wide

Sentence boundaries ✓   Coreference resolution ✓   Context ✓

Semantic relationships ✓   Part of speech ✓   Comparisons ✓

80

# Feed Forward

### Feed Forward

- Non Linearity

- Complex Relationships

- Learn from each other

**Feed Forward**

Residuals

Input                Norm(Z)

# Add & Norm



Input

$+$

Norm(Z)

**Normalization**

Mean 0, Std dev 1

Stabilizes training

Regularization effect

**Add Residuals**

Avoid vanishing gradients

Train deeper networks

N×

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

# Add & Norm



Add & Norm

Feed Forward

Input          Residuals          Norm(Z)

# Encoders

**Encoder**

ENCODER

# Encoders

**Encoder**

ENCODER

.
.
.

ENCODER

ENCODER

Input to Encoder$_{i+1}$

↑

Output from Encoder$_i$



Add & Norm

Feed Forward

N×

Add & Norm

Multi-Head Attention

Positional Encoding ⊕

Input Embedding

Inputs

# Transformers

✔ **Tokenization**

✔ **Input Embeddings**

✔ **Position Encodings**

✔ **Query, Key, & Value**

✔ **Attention**

✔ **Self Attention**

✔ **Multi-Head Attention**

✔ **Feed Forward**

✔ **Add & Norm**

✔ **Encoders**

- Masked Attention
- Encoder Decoder Attention
- Linear
- Softmax
- Decoders
- Encoder-Decoder Models

# Machine Translation

# Targets

Targets

Ich habe einen Apfel gegessen

# Targets



**Embedding Layer + Positional Encoding**

| <sos> | Ich | habe | einen | Apfel | gegessen | <eos> |

**Tokenizer**

Ich habe einen Apfel gegessen

Generate Target Embeddings

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Feed Forward

N×

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

N×

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

# Masked Multi Head Attention



| <sos> | Ich | habe | einen | Apfel | gegessen | <eos> |

# Masked Multi Head Attention

Decoding step by step (using Teacher Forcing)

Inference

# Masked Multi Head Attention

Decoding step by step (using Teacher Forcing)



Inference

1  <sos>
2  <sos>  Ich
3  <sos>  Ich
4  <sos>
5  <sos>
6  <sos>  Apfel  gegessen
7  <sos>  Ic  habe  einen  Apfel  gegessen  <eos>

**Parallelized ?**

92

# Masked Multi Head Attention

Decoding step by step (using Teacher Forcing)

Training

| <sos> | Ich | habe | einen | Apfel | gegessen | <eos> |

# Masked Multi Head Attention

Decoding step by step (using Teacher Forcing)

Training

| <sos> | Ich | habe | einen | Apfel | gegessen | <eos> |

*Outputs at time T should only pay attention to outputs*

*until time T-1*

# Masked Multi Head Attention

Decoding step by step (using Teacher Forcing)



| | | | | | | |
|---|---|---|---|---|---|---|
| 1 | <sos> | Ich | habe | einen | Apfel | gegessen | <eos> |
| 2 | <sos> | Ich | habe | einen | Apfel | gegessen | <eos> |
| 3 | <sos> | Ich | habe | einen | Apfel | gegessen | <eos> |
| 4 | <sos> | Ich | habe | einen | Apfel | gegessen | <eos> |
| 5 | <sos> | Ich | habe | einen | Apfel | gegessen | <eos> |
| 6 | <sos> | Ich | habe | einen | Apfel | gegessen | <eos> |
| 7 | <sos> | Ich | habe | einen | Apfel | gegessen | <eos> |

# Masked Multi Head Attention

Decoding step by step (using Teacher Forcing)

| | | | | | | |
|---|---|---|---|---|---|---|
| 1 | <sos> | Ich | habe | einen | Apfel | gegessen | <eos> |
| 2 | <sos> | Ich | habe | einen | Apfel | gegessen | <eos> |
| 3 | <sos> | Ich | habe | einen | Apfel | gegessen | <eos> |
| 4 | <sos> | Ich | habe | einen | Apfel | gegessen | <eos> |
| 5 | <sos> | Ich | habe | einen | Apfel | gegessen | <eos> |
| 6 | <sos> | Ich | habe | einen | Apfel | gegessen | <eos> |
| 7 | <sos> | Ich | habe | einen | Apfel | gegessen | <eos> |

Mask the available attention values ?

# Masked Multi Head Attention

Decoding step by step (using Teacher Forcing)



| | | | | | | |
|---|---|---|---|---|---|---|
| 1 | <sos> | - ∞ | - ∞ | - ∞ | - ∞ | - ∞ |
| 2 | <sos> | Ich | - ∞ | - ∞ | - ∞ | - ∞ |
| 3 | <sos> | Ich | habe | - ∞ | - ∞ | - ∞ |
| 4 | <sos> | Ich | habe | einen | - ∞ | - ∞ |
| 5 | <sos> | Ich | habe | einen | Apfel | - ∞ |
| 6 | <sos> | Ich | habe | einen | Apfel | gegessen |
| 7 | <sos> | Ich | habe | einen | Apfel | gegessen |

# Masked Multi Head Attention

Decoding step by step (using Teacher Forcing)



| | | | | | | |
|---|---|---|---|---|---|---|
| 1 | <sos> | - ∞ | - ∞ | - ∞ | - ∞ | - ∞ | - ∞ |
| 2 | <sos> | Ich | - ∞ | - ∞ | - ∞ | - ∞ | - ∞ |
| 3 | <sos> | Ich | habe | - ∞ | - ∞ | - ∞ | - ∞ |
| 4 | <sos> | Ich | habe | einen | - ∞ | - ∞ | - ∞ |
| 5 | <sos> | Ich | habe | einen | Apfel | - ∞ | - ∞ |
| 6 | <sos> | Ich | habe | einen | Apfel | gegessen | - ∞ |
| 7 | <sos> | Ich | habe | einen | Apfel | gegessen | <eos> |

Softmax    - ∞    -> 0

# Masked Multi Head Attention

$R^{T \times T}$     $R^{T \times T}$

$QK^T$     Attention Mask: M     Masked Attention

## Masked Multi Head Attention : Z'

# Masked Multi Head Attention

$R^{T \times T}$

$R^{T \times d_h}$

$\otimes$

Masked Attention

Values

Masked Multi Head Attention : Z'

Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

N×

Add & Norm

Feed
Forward

Add & Norm

Masked
Multi-Head
Attention

N×

Add & Norm

Multi-Head
Attention

Positional
Encoding

Positional
Encoding

Input
Embedding

Output
Embedding

Inputs

Outputs
(shifted right)

# Encoder Decoder Attention



Encoder Decoder Attention ?

Add & Norm

Input

Norm(Z')

# Encoder Decoder Attention

**Encoder Decoder Attention ?**

# Encoder Decoder Attention

Encoder **Self** Attention

1. Queries from Encoder Inputs
2. Keys from Encoder Inputs
3. Values from Encoder Inputs

Decoder **Masked Self** Attention

1. Queries from Decoder Inputs
2. Keys from Decoder Inputs
3. Values from Decoder Inputs

# Attention

{Key, Value store}

{Query: "Order details of order_104"}

{Query: "Order details of order_106"}

{"order_100": {"items":"a1", "delivery_date":"a2", ....}},
{"order_101": {"items":"b1", "delivery_date":"b2", ...}},
{"order_102": {"items":"c1", "delivery_date":"c2", ...}},
{"order_103": {"items":"d1", "delivery_date":"d2", ...}},
{"order_104": {"items":"e1", "delivery_date":"e2", ...}},
{"order_105": {"items":"f1", "delivery_date":"f2", ...}},
{"order_106": {"items":"g1", "delivery_date":"g2", ...}},
{"order_107": {"items":"h1", "delivery_date":"h2", ...}},
{"order_108": {"items":"i1", "delivery_date":"i2", ...}},
{"order_109": {"items":"j1", "delivery_date":"j2", ...}},
{"order_110": {"items":"k1", "delivery_date":"k2", ....}}

# Encoder Decoder Attention

**Encoder**

Keys from **Encoder Outputs**
Values from **Encoder Outputs**

**Decoder**

Queries from **Decoder Inputs**

NOTE: Every decoder block receives the same FINAL encoder output

# Encoder Decoder Attention

- Non Linearity

- Complex Relationships

- Learn from each other



Feed Forward

Residuals

Norm(Z'')

Add n Norm Decoder Self Attn

# Decoder

DECODER



Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

N×

Add & Norm

Masked
Multi-Head
Attention

Positional
Encoding

Output
Embedding

Outputs
(shifted right)

Add & Norm

Feed
Forward

N×

Add & Norm

Multi-Head
Attention

Positional
Encoding

Input
Embedding

Inputs

# Decoder

**DECODER**

.
.
.

**DECODER**

**DECODER**

$R^{T_d \times d_{model}}$



Decoder output

# Linear

$R^{V \times d_{model}}$

**Linear weights are often tied with input embedding matrix**



$R^{T_d \times d_{model}}$

softmax

$\bigotimes$

...

Final Decoder Output

$R^{T_d \times V}$

Linear

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Add & Norm

Feed Forward

N×

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

# Softmax

Output Probabilities

V $\longrightarrow$

$T_d$

$R^{T_d \times V}$

# Poll 2 - @1580

**Which of the following are true about transformers?**

a. Transformers can always be run in parallel

b. Transformer decoders can only be parallelized during training

c. Queries, keys, and values are obtained by splitting the input into 3 equal segments

d. Multihead attention might help transformers find different kinds of relations between tokens

e. Decoder outputs provide attention queries and keys, while the values come from the encoder

# Poll 2 - @1580

**Which of the following are true about transformers?**

   a.     Transformers can always be run in parallel

   b.     **Transformer decoders can only be parallelized during training**

   c.     Queries, keys, and values are obtained by splitting the input into 3 equal segments

   d.     **Multihead attention might help transformers find different kinds of relations between tokens**

   e.     Decoder outputs provide attention queries and keys, while the values come from the encoder

# Transformers

Targets

Ich habe einen Apfel gegessen

Inputs

I ate an apple

**Machine Translation**

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Add & Norm

Masked Multi-Head Attention

Add & Norm

Feed Forward

N×

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

Positional Encoding

Output Embedding

Outputs (shifted right)

# Transformers

✔ **Tokenization**

✔ **Input Embeddings**

✔ **Position Encodings**

✔ **Query, Key, & Value**

✔ **Attention**

✔ **Self Attention**

✔ **Multi-Head Attention**

✔ **Feed Forward**

✔ **Add & Norm**

✔ **Encoders**

✔ **Masked Attention**

✔ **Encoder Decoder Attention**

✔ **Linear**

✔ **Softmax**

✔ **Decoders**

• Encoder-Decoder Models

# Transformers

# Transformers



**Representation**

**Generation**

116

# Transformers



**Input** – input tokens
**Output** – hidden states

**Representation**

**Input** – output tokens and hidden states*
**Output** – output tokens

**Generation**

# Transformers



**Input** – input tokens
**Output** – hidden states

**Model can see all timesteps**

**Representation**

**Input** – output tokens and hidden states*
**Output** – output tokens

**Model can only see previous timesteps**

**Generation**

118

# Transformers



**Input** – input tokens
**Output** – hidden states

**Model can see all timesteps**

**Does not usually output tokens, so no inherent auto-regressivity**

**Representation**

**Input** – output tokens and hidden states*
**Output** – output tokens

**Model can only see previous timesteps**

**Model is auto-regressive with previous timesteps' outputs**

**Generation**

# Transformers

**Input** – input tokens
**Output** – hidden states

**Model can see all timesteps**

**Does not usually output tokens, so no inherent auto-regressivity**

*Can also be adapted to generate tokens by appending a module that maps hidden state dimensionality to vocab size*

Representation



**Input** – output tokens and hidden states*
**Output** – output tokens

**Model can only see previous timesteps**

**Model is auto-regressive with previous timesteps' outputs**

*Can also be adapted to generate hidden states by looking before token outputs*

Generation

# Transformers

✔ **Tokenization**

✔ **Input Embeddings**

✔ **Position Encodings**

✔ **Query, Key, & Value**

✔ **Attention**

✔ **Self Attention**

✔ **Multi-Head Attention**

✔ **Feed Forward**

✔ **Add & Norm**

✔ **Encoders**

✔ **Masked Attention**

✔ **Encoder Decoder Attention**

✔ **Linear**

✔ **Softmax**

✔ **Decoders**

✔ **Encoder-Decoder Models**

# Part 2

# Pre-training and Fine-tuning

# How to train and fine-tune transformers

**1. Training**

**2. Inference**

# How to train and fine-tune transformers

**1. Pre-training** → **2. Fine-tuning** → **3. Inference**

# How to train and fine-tune transformers



**1. Pre-training**

Transformer architecture

Supervised training

Larger general dataset

Pre-trained transformer model

**2. Fine-tuning**

**3. Inference**

**Lot's of data, learn general things. May serve as a parameter initialization.**

**Usually requires significant computational resources and time.**

# How to train and fine-tune transformers



**1. Pre-training**

Transformer architecture

Supervised training

Pre-trained transformer model

**2. Fine-tuning**

Task-specific training

Fine-tuned transformer model

**3. Inference**

Larger general dataset

Smaller task-specific dataset

**Lot's of data, learn general things. May serve as a parameter initialization.**

**Usually requires significant computational resources and time.**

**Adaptation to the specific task.**

**Potentially less computationally intensive.**

# Parameter-Efficient Fine-Tuning Techniques



LoRA: https://arxiv.org/abs/2106.09685
BitFit: https://arxiv.org/abs/2106.10199

# Parameter-Efficient Fine-Tuning Techniques

**LoRA (Lower-Rank Adaptation)**



Figure 1: Our reparametrization. We only train $A$ and $B$.



LoRA: https://arxiv.org/abs/2106.09685
BitFit: https://arxiv.org/abs/2106.10199

# Parameter-Efficient Fine-Tuning Techniques

## LoRA (Lower-Rank Adaptation)



Figure 1: Our reparametrization. We only train $A$ and $B$.

Pretrained Weights $W \in \mathbb{R}^{d \times d}$

$B = 0$

$A = \mathcal{N}(0, \sigma^2)$

## BitFit

$$\mathbf{Q}^{m,\ell}(\mathbf{x}) = \mathbf{W}_q^{m,\ell}\mathbf{x} + \mathbf{b}_q^{m,\ell}$$

$$\mathbf{K}^{m,\ell}(\mathbf{x}) = \mathbf{W}_k^{m,\ell}\mathbf{x} + \mathbf{b}_k^{m,\ell}$$

$$\mathbf{V}^{m,\ell}(\mathbf{x}) = \mathbf{W}_v^{m,\ell}\mathbf{x} + \mathbf{b}_v^{m,\ell}$$



LoRA: https://arxiv.org/abs/2106.09685
BitFit: https://arxiv.org/abs/2106.10199

# Part 3

## Transformer Applications

# Transformers



**Representation / Encoder**

**Generation / Decoder**

# Data Modalities

- **Language (see Part 4 of the lecture)**

- **Vision**

- **Audio**

- **… and many other modalities (e.g., biological/physiological signals, etc.)**

- **Multimodal (>2 data modalities)**

# Computer Vision

1. In computer vision convolutional architectures remain largely dominant.

2. Inspired by NLP successes, multiple works try introducing combining CNN-like architectures with self-attention or replacing the convolutions entirely.

3. However, they faced challenges with performance and scaling.

4. Key breakthrough - Vision Transformer (ViT) released in 2020

An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale

# Computer Vision - Tokenization



An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale

# Vision Transformer (ViT) Model Architecture



**Vision Transformer (ViT)**

Class
Bird
Ball
Car
...

MLP Head

Transformer Encoder

Patch + Position Embedding

* Extra learnable [class] embedding

Linear Projection of Flattened Patches

1. Split an image into fixed-size patches (16x16 pixels).

2. Tokenize each path (linear projection of flattened patches).

3. Add position embedding.

4. Feed the resulting sequence of vectors to a standard Transformer encoder.

5. For classification, add an extra learnable"classification token" to the sequence.

**An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale**

# ViT - Learning Patterns



- ViT learns the grid like structure of the image patches via its position embeddings.
- The lower layers contain both global and local features, the higher layers contain only global features.

Transformers for Image Recognition at Scale

# ViT Performance

| | ViT-H | Previous SOTA |
|---|---|---|
| **ImageNet** | 88.55 | 88.5 |
| **ImageNet-ReaL** | 90.72 | 90.55 |
| **Cifar-10** | 99.50 | 99.37 |
| **Cifar-100** | 94.55 | 93.51 |
| **Pets** | 97.56 | 96.62 |
| **Flowers** | 99.68 | 99.63 |



- ViT model attains state-of-the-art performance on multiple popular benchmarks, including 88.55% top-1 accuracy on ImageNet and 99.50% on CIFAR-10

Transformers for Image Recognition at Scale

# Audio

- Similar to the computer vision but with spectrograms instead of images.

- Exists as encoder-decoder variants or as an encoder-only variant with CTC loss.

- Could be augmented with the CNN.

Conformer: Convolution-augmented Transformer for Speech Recognition

AST: Audio Spectrogram Transformer

# Multimodal Transformer - UniT



UniT: Multimodal Multitask Learning with a Unified Transformer

1. UniT handles 7 tasks ranging from object detection to vision-and language reasoning and natural language understanding.

2. Components:

   ○ An image encoder to encode the visual inputs.

   ○ A text encoder to encode the language inputs.

   ○ A joint decoder with per-task query embedding.

   ○ Task-specific heads to make the final outputs for each task.

# Multimodal Transformer - LLaVA



Language Response $X_a$

Language Model $f_\phi$

Projection $W$

$Z_v$

$H_v$

Vision Encoder

$X_v$ Image

$H_q$

$X_q$ Language Instruction

**Visual Instruction Tuning (LLaVA - Large Language and Vision Assistant)**

# Multimodal Transformer - LLaVA



Start a new conversation, and the history is cleared.

**User:** Do you know who drew this painting?

**LLaVA:** The painting depicts a dog in a humorous situation, where it appears to be wearing a woman's dress and standing in front of a landscape. Although the artist of the painting is not explicitly identified, it could be a reference to the famous artist Da Vinci, known for his artistic works, including the Mona Lisa. The painting is likely a creative and comical take on the traditional portrait style, replacing a human figure with a dog dressed in a woman's clothing, making it an interesting and amusing visual piece.

Visual Instruction Tuning (LLaVA - Large Language and Vision Assistant)

# Part 4

# Large Language Models

# Transformers, mid-2017



**Representation**

**Generation**

# 2018 – Inception of the LLM Era



**BERT
Oct 2018**

**Representation**

**GPT
June 2018**

**Generation**

# BERT - Bidirectional Encoder Representations

- One of the biggest challenges in LM-building used to be the lack of task-specific training data.

- What if we learn an effective representation that can be applied to a variety of downstream tasks?
  - Word2vec (2013)
  - GloVe (2014)

# BERT - Bidirectional Encoder Representations

**BERT Pre-Training Corpus:**

- English Wikipedia - 2,500 million words
- Book Corpus - 800 million words

# BERT - Bidirectional Encoder <mark>Representations</mark>

**BERT Pre-Training Corpus:**

- English Wikipedia - 2,500 million words
- Book Corpus - 800 million words

**BERT Pre-Training Tasks:**

- MLM (Masked Language Modeling)
- NSP (Next Sentence Prediction)

# BERT - Bidirectional Encoder Representations

**BERT Pre-Training Corpus:**
- English Wikipedia - 2,500 million words
- Book Corpus - 800 million words

**BERT Pre-Training Tasks:**
- MLM (Masked Language Modeling)
- NSP (Next Sentence Prediction)

**BERT Pre-Training Results:**
- BERT-Base – 110M Params
- BERT-Large – 340M Params

# BERT - Bidirectional Encoder Representations

## MLM (Masked Language Modeling)

# BERT - Bidirectional Encoder **Representations**

**NSP (Next Sentence Prediction)**

# BERT - Bidirectional Encoder <mark>Representations</mark>

**BERT Fine-Tuning:**

- Simply add a task-specific module after the last encoder layer to map it to the desired dimension.

  - *Classification Tasks:*
    - *Add a feed-forward layer on top of the encoder output for the [CLS] token*
  - *Question Answering Tasks:*
    - *Train two extra vectors to mark the beginning and end of answer from paragraph*

  - *…*

# BERT - Bidirectional Encoder Representations

## BERT Evaluation:

- General Language Understanding Evaluation (GLUE)
  - Sentence pair tasks
  - Single sentence classification

- Stanford Question Answering Dataset (SQuAD)

# BERT - Bidirectional Encoder Representations

## BERT Evaluation:

| System | MNLI-(m/mm) 392k | QQP 363k | QNLI 108k | SST-2 67k | CoLA 8.5k | STS-B 5.7k | MRPC 3.5k | RTE 2.5k | Average - |
|---|---|---|---|---|---|---|---|---|---|
| Pre-OpenAI SOTA | 80.6/80.1 | 66.1 | 82.3 | 93.2 | 35.0 | 81.0 | 86.0 | 61.7 | 74.0 |
| BiLSTM+ELMo+Attn | 76.4/76.1 | 64.8 | 79.8 | 90.4 | 36.0 | 73.3 | 84.9 | 56.8 | 71.0 |
| OpenAI GPT | 82.1/81.4 | 70.3 | 87.4 | 91.3 | 45.4 | 80.0 | 82.3 | 56.0 | 75.1 |
| BERT$_{BASE}$ | 84.6/83.4 | 71.2 | 90.5 | 93.5 | 52.1 | 85.8 | 88.9 | 66.4 | 79.6 |
| BERT$_{LARGE}$ | **86.7/85.9** | **72.1** | **92.7** | **94.9** | **60.5** | **86.5** | **89.3** | **70.1** | **82.1** |

| System | Dev EM | Dev F1 | Test EM | Test F1 |
|---|---|---|---|---|
| Leaderboard (Oct 8th, 2018) | | | | |
| Human | - | - | 82.3 | 91.2 |
| #1 Ensemble - nlnet | - | - | 86.0 | 91.7 |
| #2 Ensemble - QANet | - | - | 84.5 | 90.5 |
| #1 Single - nlnet | - | - | 83.5 | 90.1 |
| #2 Single - QANet | - | - | 82.5 | 89.3 |
| Published | | | | |
| BiDAF+ELMo (Single) | - | 85.8 | - | - |
| R.M. Reader (Single) | 78.9 | 86.3 | 79.5 | 86.6 |
| R.M. Reader (Ensemble) | 81.2 | 87.9 | 82.3 | 88.5 |
| Ours | | | | |
| BERT$_{BASE}$ (Single) | 80.8 | 88.5 | - | - |
| BERT$_{LARGE}$ (Single) | 84.1 | 90.9 | - | - |
| BERT$_{LARGE}$ (Ensemble) | 85.8 | 91.8 | - | - |
| BERT$_{LARGE}$ (Sgl.+TriviaQA) | **84.2** | **91.1** | **85.1** | **91.8** |
| BERT$_{LARGE}$ (Ens.+TriviaQA) | **86.2** | **92.2** | **87.4** | **93.2** |

Table 2: SQuAD results. The BERT ensemble is 7x systems which use different pre-training checkpoints and fine-tuning seeds.



153

# BERT - Bidirectional Encoder <mark>Representations</mark>

**What is our takeaway from BERT?**

- **Pre-training tasks can be invented flexibly…**
  - Effective representations can be derived from a flexible regime of pre-training tasks.

# BERT - Bidirectional Encoder <mark>Representations</mark>

**What is our takeaway from BERT?**

- **Pre-training tasks can be invented flexibly…**
  - Effective representations can be derived from a flexible regime of pre-training tasks.

- **Different NLP tasks seem to be highly transferable with each other...**
  - As long as we have effective representations, that seems to form a general model which can serve as the backbone for many specialized models.

# BERT - Bidirectional Encoder <mark>Representations</mark>

**What is our takeaway from BERT?**

- **Pre-training tasks can be invented flexibly…**
  - Effective representations can be derived from a flexible regime of pre-training tasks.

- **Different NLP tasks seem to be highly transferable with each other...**
  - As long as we have effective representations, that seems to form a general model which can serve as the backbone for many specialized models.

- **And scaling works!!!**
  - 340M was considered large in 2018

# 2018 – Inception of the LLM Era



**BERT
Oct 2018**

**Representation**

**GPT
June 2018**

**Generation**

# GPT – ==Generative== Pretrained Transformer

- Similarly motivated as BERT, though differently designed

  - Can we leverage large amounts of unlabeled data to pretrain an LM that understands general patterns?

# GPT – ==Generative== Pretrained Transformer

**GPT Pre-Training Corpus:**
- Similarly, BooksCorpus and English Wikipedia

**GPT Pre-Training Tasks:**
- Predict the next token, given the previous tokens
  - More learning signals than MLM

**GPT Pre-Training Results:**
- GPT – 117M Params
  - Similarly competitive on GLUE and SQuAD

# GPT – <mark>Generative</mark> Pretrained Transformer

## GPT Fine-Tuning:

- Prompt-format task-specific text as a continuous stream for the model to fit

**Summarization**

Summarize this article:

The summary is:

**QA**

Answer the question based on the context.

Context:

Question:

Answer:

# GPT – <mark>Generative</mark> Pretrained Transformer

**What is our takeaway from GPT?**

- **The Effectiveness of Self-Supervised Learning**
  - Specifically, the model seems to be able to learn from generating the language *itself*, rather than from any specific task we might cook up.

# GPT – <mark>Generative</mark> Pretrained Transformer

**What is our takeaway from GPT?**

- **The Effectiveness of Self-Supervised Learning**
  - Specifically, the model seems to be able to learn from generating the language *itself*, rather than from any specific task we might cook up.

- **Language Model as a Knowledge Base**
  - Specifically, a generatively pretrained model seems to have a decent zero-shot performance on a range of NLP tasks.

# GPT – ==Generative== Pretrained Transformer

**What is our takeaway from GPT?**

- **The Effectiveness of Self-Supervised Learning**
  - Specifically, the model seems to be able to learn from generating the language *itself*, rather than from any specific task we might cook up.

- **Language Model as a Knowledge Base**
  - Specifically, a generatively pretrained model seems to have a decent zero-shot performance on a range of NLP tasks.

- **And scaling works!!!**



163

# Poll 3 - @1579

**The original GPT's parameter count is closest to…**

A. 117
B. 117K
C. 117M
D. 117B

# Poll 3 - @1579

**The original GPT's parameter count is closest to…**

A. 117
B. 117K
**C. 117M**
D. 117B

# The LLM Era – Paradigm Shift in Machine Learning



BERT
Oct 2018

Representation

GPT
June 2018

Generation

# The LLM Era – Paradigm Shift in Machine Learning



**BERT** – 2018
**DistilBERT** – 2019
**RoBERTa** – 2019
**ALBERT** – 2019
**ELECTRA** – 2020
**DeBERTa** – 2020

…

**Representation**

**T5** – 2019
**BART** – 2019
**mT5** – 2021

…

**GPT** – 2018
**GPT-2** – 2019
**GPT-3** – 2020
**GPT-Neo** – 2021
**GPT-3.5 (ChatGPT)** – 2022
**LLaMA** – 2023
**GPT-4** – 2023

…

**Generation**

# The LLM Era – Paradigm Shift in Machine Learning

**From both BERT and GPT, we learn that…**
- Transformers seem to provide a new class of generalist models that are capable of capturing knowledge which is more fundamental than task-specific abilities.

<u>**Before LLMs**</u>                                                    <u>**Since LLMs**</u>

- **Feature Engineering**
  - How do we design or select the best features for a task?

# The LLM Era – Paradigm Shift in Machine Learning

**From both BERT and GPT, we learn that…**
- Transformers seem to provide a new class of generalist models that are capable of capturing knowledge which is more fundamental than task-specific abilities.

<u>**Before LLMs**</u>                                          <u>**Since LLMs**</u>

- **Feature Engineering**
  - How do we design or select the best features for a task?
- <span style="color:darkred">**Model Selection**</span>
  - <span style="color:darkred">Which model is best for which type of task?</span>

# The LLM Era – Paradigm Shift in Machine Learning

**From both BERT and GPT, we learn that…**
- Transformers seem to provide a new class of generalist models that are capable of capturing knowledge which is more fundamental than task-specific abilities.

|  Before LLMs | Since LLMs |
|---|---|

- **Feature Engineering**
  - How do we design or select the best features for a task?
- **Model Selection**
  - Which model is best for which type of task?
- **Transfer Learning**
  - Given scarce labeled data, how do we transfer knowledge from other domains?

# The LLM Era – Paradigm Shift in Machine Learning

**From both BERT and GPT, we learn that…**

- Transformers seem to provide a new class of generalist models that are capable of capturing knowledge which is more fundamental than task-specific abilities.

<u>**Before LLMs**</u>                                                <u>**Since LLMs**</u>

- **Feature Engineering**
  - How do we design or select the best features for a task?
- **Model Selection**
  - Which model is best for which type of task?
- **Transfer Learning**
  - Given scarce labeled data, how do we transfer knowledge from other domains?
- **Overfitting vs Generalization**
  - How do we balance complexity and capacity to prevent overfitting while maintaining good performance?

# The LLM Era – Paradigm Shift in Machine Learning

**From both BERT and GPT, we learn that…**

- Transformers seem to provide a new class of generalist models that are capable of capturing knowledge which is more fundamental than task-specific abilities.

<table>
<tr><th>Before LLMs</th><th>Since LLMs</th></tr>
<tr><td>

- **Feature Engineering**
  - How do we design or select the best features for a task?
- **Model Selection**
  - Which model is best for which type of task?
- **Transfer Learning**
  - Given scarce labeled data, how do we transfer knowledge from other domains?
- **Overfitting vs Generalization**
  - How do we balance complexity and capacity to prevent overfitting while maintaining good performance?

</td><td>

- **Pre-training and Fine-tuning**
  - How do we leverage large scales of unlabeled data out there previously under-leveraged?

</td></tr>
</table>

# The LLM Era – Paradigm Shift in Machine Learning

**From both BERT and GPT, we learn that…**

- Transformers seem to provide a new class of generalist models that are capable of capturing knowledge which is more fundamental than task-specific abilities.

<div>

**Before LLMs**

- **Feature Engineering**
  - How do we design or select the best features for a task?
- **Model Selection**
  - Which model is best for which type of task?
- **Transfer Learning**
  - Given scarce labeled data, how do we transfer knowledge from other domains?
- **Overfitting vs Generalization**
  - How do we balance complexity and capacity to prevent overfitting while maintaining good performance?

</div>

<div>

**Since LLMs**

- **Pre-training and Fine-tuning**
  - How do we leverage large scales of unlabeled data out there previously under-leveraged?
- **Zero-shot and Few-shot learning**
  - How can we make models perform on tasks they are _not_ trained on?

</div>

# The LLM Era – Paradigm Shift in Machine Learning

**From both BERT and GPT, we learn that…**
- Transformers seem to provide a new class of generalist models that are capable of capturing knowledge which is more fundamental than task-specific abilities.

<table>
<tr><td align="center"><u>Before LLMs</u></td><td align="center"><u>Since LLMs</u></td></tr>
<tr><td valign="top">

- **Feature Engineering**
  - How do we design or select the best features for a task?
- **Model Selection**
  - Which model is best for which type of task?
- **Transfer Learning**
  - Given scarce labeled data, how do we transfer knowledge from other domains?
- **Overfitting vs Generalization**
  - How do we balance complexity and capacity to prevent overfitting while maintaining good performance?

</td><td valign="top">

- **Pre-training and Fine-tuning**
  - How do we leverage large scales of unlabeled data out there previously under-leveraged?
- **Zero-shot and Few-shot learning**
  - How can we make models perform on tasks they are <u>*not*</u> trained on?
- **Prompting**
  - How do we make models understand their task simply by describing it in natural language?

</td></tr>
</table>

# The LLM Era – Paradigm Shift in Machine Learning

**From both BERT and GPT, we learn that…**

- Transformers seem to provide a new class of generalist models that are capable of capturing knowledge which is more fundamental than task-specific abilities.

<u>Before LLMs</u>

- **Feature Engineering**
  - How do we design or select the best features for a task?
- **Model Selection**
  - Which model is best for which type of task?
- **Transfer Learning**
  - Given scarce labeled data, how do we transfer knowledge from other domains?
- **Overfitting vs Generalization**
  - How do we balance complexity and capacity to prevent overfitting while maintaining good performance?

<u>Since LLMs</u>

- **Pre-training and Fine-tuning**
  - How do we leverage large scales of unlabeled data out there previously under-leveraged?
- **Zero-shot and Few-shot learning**
  - How can we make models perform on tasks they are <u>not</u> trained on?
- **Prompting**
  - How do we make models understand their task simply by describing it in natural language?
- **Interpretability and Explainability**
  - How can we <u>understand</u> the inner workings of our own models?

# The LLM Era – Paradigm Shift in Machine Learning

- **What has caused this paradigm shift?**

# The LLM Era – Paradigm Shift in Machine Learning

- **What has caused this paradigm shift?**

  - **Recall: Problem in recurrent networks**
    - Information is effectively lost during encoding of long sequences
    - Sequential nature disables parallel training and favors late timestep inputs

# The LLM Era – Paradigm Shift in Machine Learning

- **What has caused this paradigm shift?**

  - **Recall: Problem in recurrent networks**
    - Information is effectively lost during encoding of long sequences
    - Sequential nature disables parallel training and favors late timestep inputs

  - **Solution: Attention is all you need!!!**
    - Handling long-range dependencies
    - Parallel training
    - Dynamic attention weights based on inputs

# The LLM Era – Paradigm Shift in Machine Learning

- **Attention and Transformer – is this the end?**

# The LLM Era – Paradigm Shift in Machine Learning

- **Attention and Transformer – is this the end?**

  - **Problem in current Transformer-based LLMs??**

# The LLM Era – Paradigm Shift in Machine Learning

- **Attention and Transformer – is this the end?**

  - **Problem in current Transformer-based LLMs??**
    - True understanding the material vs. memorization and pattern-matching
    - Cannot reliably follow rules – factual hallucination e.g. inability in arithmetic

# The LLM Era – Paradigm Shift in Machine Learning

- **Attention and Transformer – is this the end?**

  - **Problem in current Transformer-based LLMs??**
    - True understanding the material vs. memorization and pattern-matching
    - Cannot reliably follow rules – factual hallucination e.g. inability in arithmetic

  - **Solution: ???**

# Looking Back

It is true that language models are just programmed to predict the next token…

In fact, all animals, including us, are just programmed to survive and reproduce, and yet amazingly complex and beautiful stuff comes from it.

**- Sam Altman***

**\*Paraphrased**