

# Transformer and Large-Language Models

**Tianxiang (Adam) Gao**

November 4, 2024

# Outline

- ## 2 Large-Language Models

## Recap: Seq2Seq Models

- **Word Analogy Task:** Given an analogy “a is to b as c is to ?”, find the word  $d$  by maximizing cosine similarity.
- **GloVe:** Uses a co-occurrence matrix to capture word meaning based on surrounding context.
- **Bias in Word Embeddings:** Identify bias directions using definitional pairs, then neutralize non-definitional words.
- **Seq2Seq:** Use an RNN Encoder-Decoder architecture to handle tasks where both input and output are sequences.
- **Conditional Language Model:** The output sequence is generated sequentially based on the context vector that summarizes the input sequence
- **Beam Search:** Keeps multiple high-probability sequences to improve output quality.
- **BLEU Score:** A metric using modified precision to assess the accuracy of generated sequences.
- **Distinct Convex Vector:** A distinct context word  $e_t$  is used to generate each target word  $\hat{y}_t$

$$\mathbb{P}(y_t \mid x, y_1, \dots, y_{t-1}) = \mathbb{P}_\phi(y_t \mid s_t), \quad \text{where} \quad s_t = g_\phi(s_{t-1}, y_{t-1}, c_t)$$

- **Attention Weights:** The distinct context word  $e_t$  is a weighted sum of encoder hidden states:

$$c_t = \sum_i \alpha_{t,i} h_i,$$

where  $\alpha_{t,i} = \text{softmax}(e_{t,i})$  are **attention weights**

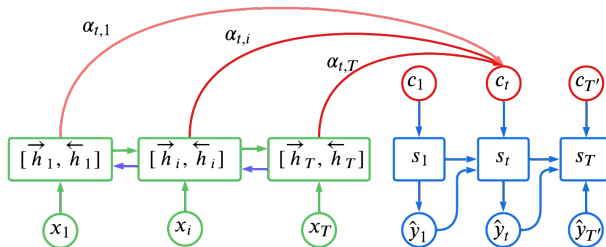
- **Alignment Scores:** Alignment scores  $e_{t,i} = a(s_{t-1}, h_i)$  indicate relevance between encoder hidden states  $h_i$  and decoder states  $s_{t-1}$ .

# Outline

1 Transformer

2 Large-Language Models

## Recall: Attention Mechanism in Seq2Seq



- In Seq2Seq models, the conditional probability of the target word  $\mathbf{y}_t$  is given by:

$$\mathbb{P}(\mathbf{y}_t \mid \mathbf{x}, \mathbf{y}_1, \dots, \mathbf{y}_{t-1}) = g_{\phi}(\mathbf{s}_t),$$

where the decoder hidden state  $\mathbf{s}_t$  is updated with the *distinct* context vector  $\mathbf{c}_t$  as follows:

$$\mathbf{s}_t = \text{GRU}(\mathbf{s}_{t-1}, [\mathbf{y}_{t-1}; \mathbf{c}_t]), \quad \text{with} \quad \mathbf{c}_t = \sum_{i=1}^{T_x} \alpha_{t,i} \mathbf{h}_i,$$

where  $\alpha_t \in \mathbb{R}^{T_x}$  are the attention weights based on alignment scores  $\mathbf{e}_t \in \mathbb{R}^{T_x}$ :

$$\alpha_t = \text{softmax}(\mathbf{e}_t), \quad \text{with} \quad \mathbf{e}_{t,i} = \mathbf{a}^{\top} \tanh(\mathbf{W}_a[\mathbf{s}_{t-1}; \mathbf{h}_i]).$$

# Self-Attention

**Define:** Self-attention creates a contextually enriched representation of each token by learning its relevance to all other tokens in the sequence.

- For each token, three vectors are computed:

$$q_t = W^q x_t, \quad k_t = W^k x_t, \quad v_t = W^v x_t$$

where the query  $q_t$  interacts with keys  $k_i$  to measure relevance:

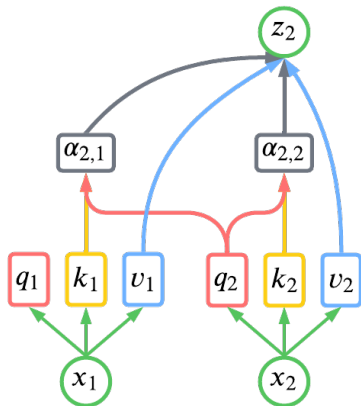
$$\alpha_{t,i} \propto q_t^\top k_i, \quad \Rightarrow \quad \alpha_t = \text{softmax} \left( \frac{K q_t}{\sqrt{d_k}} \right),$$

where

$$K = [k_1 \quad \cdots \quad k_T]^\top$$

- The new representation  $z_t$  is a weighted sum of value vectors  $v_i$ :

$$z_t = \sum_{i=1}^T \alpha_{t,i} v_i$$



# Self-Attention: Matrix Form

- Define matrices:

$$\mathbf{Q} = \mathbf{XW}^{q\top}, \quad \mathbf{K} = \mathbf{XW}^{k\top}, \quad \mathbf{V} = \mathbf{XW}^{v\top}$$

where

$$\mathbf{Q} = [\mathbf{q}_1 \quad \cdots \quad \mathbf{q}_T]^\top, \quad \mathbf{K} = [\mathbf{k}_1 \quad \cdots \quad \mathbf{k}_T]^\top, \quad \mathbf{V} = [\mathbf{v}_1 \quad \cdots \quad \mathbf{v}_T]^\top$$

- The attention weights are computed as:

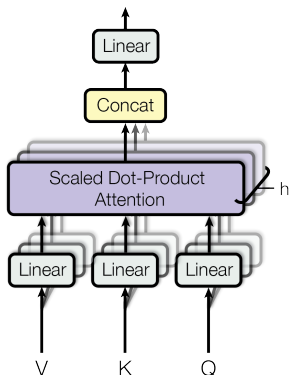
$$[\alpha_1 \quad \cdots \quad \alpha_T] = \text{softmax}\left(\frac{\mathbf{KQ}^\top}{\sqrt{d_k}}\right)$$

- The new representation  $\mathbf{Z}$  is then:

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

# Multi-Head Attention

**Definition:** Multi-head attention extends self-attention by allowing multiple heads to focus on **different aspects** of each token, capturing diverse patterns and dependencies across the sequence.



- Each head produces an independent attention output  $Z_h$ :

$$Z_h = \text{Attention}(Q_h, K_h, V_h),$$

where  $Q_h = QW_h^q$ ,  $K_h = KW_h^k$ , and  $V_h = VW_h^v$ .

- Head outputs are concatenated and linearly transformed to form the final representation:

$$Z = [Z_1 \quad \cdots \quad Z_H] W_o^T$$

where  $W_o$  is the output projection and  $H$  is the number of heads.

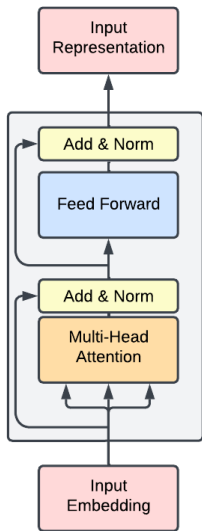
- In Matrix Form:** With  $Q_h, K_h, V_h$  for each head,

$$Z = \text{MultiHead}(Q, K, V) \in \mathbb{R}^{T \times d_{\text{model}}}$$

**Note:** The multi-head can be computed in **parallel**, each with complexity  $\mathcal{O}(T^2 d)$ .



# Multi-Head Attention Layer: LayerNorm and FNN



**Layer Normalization** computes statistics across different **hidden units**:

- In an RNN or MLP, a hidden state update is given by:

$$z = Wx, \quad h = \tanh(z)$$

where the pre-activation vector  $z \in \mathbb{R}^m$ .

- The statistics are computed across the **hidden units**:

$$z_i = w_i^\top x, \quad \mu = \frac{1}{m} \sum_{i=1}^m z_i, \quad \sigma = \sqrt{\frac{1}{m} \sum_{i=1}^m (z_i - \mu)^2}$$

where  $w_i$  is the  $i$ th row of  $W$ .

- Re-scale and shift the normalized pre-activation:

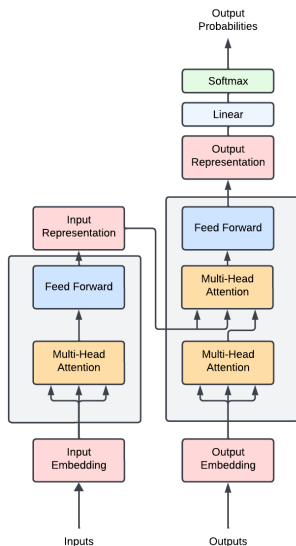
$$z_{\text{norm}} = \frac{z - \mu}{\sigma}, \quad \tilde{z} = \alpha \odot z_{\text{norm}} + \beta$$

where  $\alpha$  and  $\beta$  are trainable

**Feed-Forward Layer** captures non-linear relationships between tokens:

$$\text{FFN}(x_t) = \text{ReLU}(x_t W_1 + b_1) W_2 + b_2$$

# Multi-Head Attention Encoder and Decoder Stacks



## Encoder:

- **Input Embedding:** Converts input to dense word embeddings.
- **Multi-Head Attention:** Enhances token representations by attending to various parts of the sequence.
- **FFN:** Applies non-linear transformations to capture complex relationships between words.

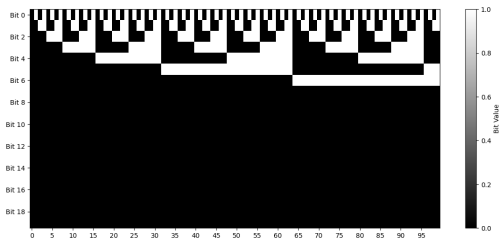
## Decoder:

- **Output Embedding:** Converts the previously generated output (one-hot encoded) into dense word embeddings.
- **First Multi-Head Attention:** Refines the output embeddings or hidden states using self-attention.
- **Second Multi-Head Attention:** Uses the output of the first attention as the query  $Q$ , with  $K$  and  $V$  from the encoder's output, allowing the decoder to attend to the input sequence.
- **Final Output:** Computes the conditional probability of the next token through a linear layer and softmax.

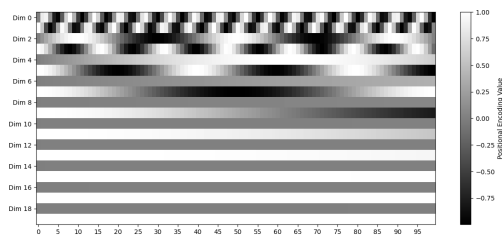
# Positional Encoding

- Unlike RNNs, transformers do **not** inherently process tokens in sequence order.
- Positional Encoding is defined as:

$$\text{PE}_{\text{pos},2i} = \sin\left(\text{pos}/10000^{\frac{2i}{d_{\text{model}}}}\right), \quad \text{and} \quad \text{PE}_{\text{pos},2i+1} = \cos\left(\text{pos}/10000^{\frac{2i}{d_{\text{model}}}}\right)$$



Binary Representation

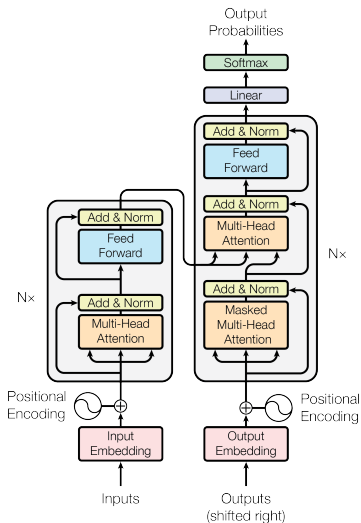


Position Encoding

## Key Properties

- Provides unique and consistent representation for each position.
- Represents positions in a low-dimensional subspace.
- Enables linear transformations for relative positioning, *i.e.*,  $\exists$  a linear  $M_k$  s.t.  $M_k \text{PE}_{\text{pos}+k} = \text{PE}_{\text{pos}}$ .

# Transformer



## Training Process: Teacher Forcing

- During training, the true output sequence is fed into the decoder.
- Masking ensures only past and current tokens are visible, preserving autoregressive properties.
- Cross-entropy loss is used to compare the predicted probability distribution with the true token.

## Loss Function: Cross-Entropy

$$\mathcal{L} = -\frac{1}{T} \sum_{t=1}^T y_t \log \mathbb{P}(\hat{y}_t)$$

- $y_t$ : True one-hot encoded token.
- $\mathbb{P}(\hat{y}_t)$ : Predicted probability for the token at time  $t$ .

# Summary

- **Self-attention** refines the representation of each token by learning its relevance to other tokens using **query-key** pairs.
- **Multi-head self-attention** captures **different aspects** of each token, enhancing the overall representation.
- **Layer normalization** computes statistics **across hidden units** to stabilize information propagation.
- **Positional encoding** adds **order** information to word embeddings, enabling the model to learn relative positioning through a linear transformation.
- **Encoder-decoder attention** refines the output representation by attending to the input sequence representation.
- The **Transformer** uses **teacher forcing** with cross-entropy loss to facilitate effective training.

# Outline

- ## 2 Large-Language Models

# BERT

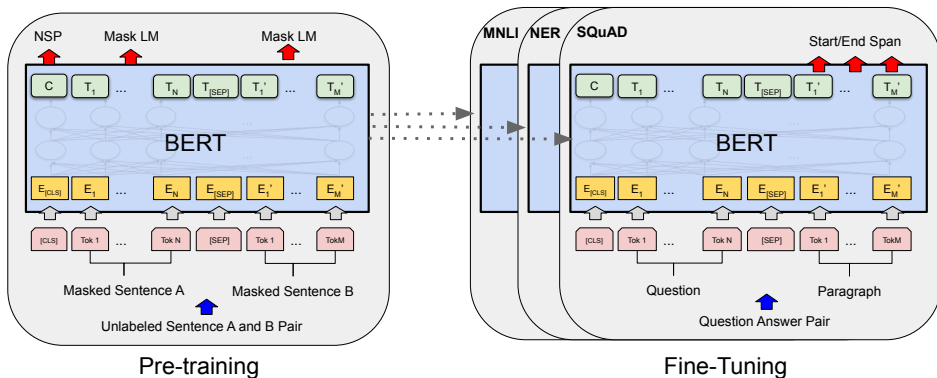


Google  
BERT

# BERT: Encoder-Only

**Definition:** BERT stands for **B**idirectional **E**ncoder **R**epresentations from **T**ransformers.

- Utilizes only the **encoder** part of the Transformer architecture.
- Designed for pre-training on large corpora and fine-tuning on downstream NLP tasks.



**Model Details:** BERT\_Large, with 340 million parameters, was trained on TPU v3 pods over 4 days using the BooksCorpus (800 million words) and English Wikipedia (2.5 billion words) datasets.



# BERT: Pre-training

**Masked Language Modeling (MLM):** Randomly masks 15% of tokens in the input sequence and predicts masked tokens based on context.

- **Problem:** [MASK] token never used in finite-tuning
- **Solution:** Do not always replace selected words with [MASK], e.g., my dog is hairy
  - 80% of the time: Replace the word with the [MASK] token, e.g., my dog is [MASK]
  - 10% of the time: Replace the word with a random word, e.g., my dog is apple
  - 10% of the time: Keep the word unchanged, e.g., my dog is hairy.

**Next Sentence Prediction (NSP):** Predicts if the second sentence follows the first, using the hidden state of the [CLS] token for binary classification (*IsNext* or *NotNext*).

- **Input**=[CLS] the man went to [MASK] store [SEP] he bought a gallon [MASK] milk [SEP], **Label**=IsNext
- **Input**=[CLS] the man [MASK] to the store [SEP] penguin [MASK] are flight ##less birds [SEP], **Label**=NotNext

## BERT: Performance on SQuAD 1.1

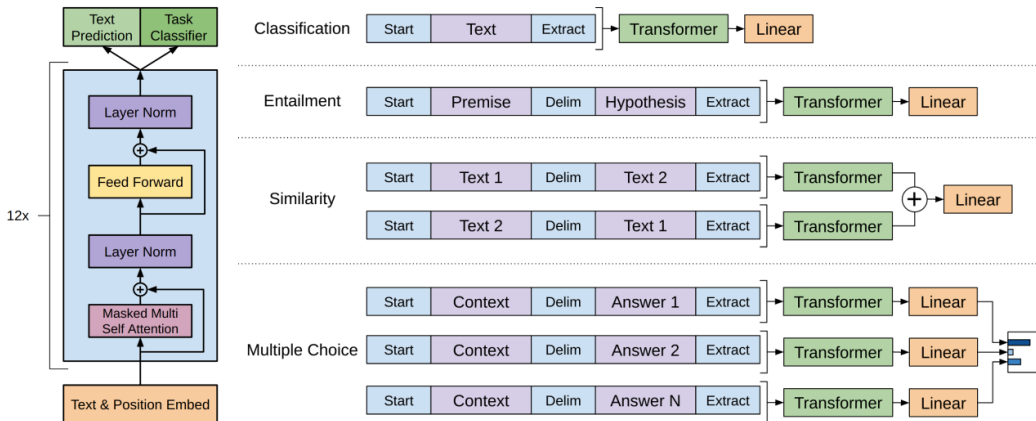
System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average -
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT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

## GPT





# GPT-1: Task-Specific Input Transformations



**Model Details:** GPT-1, with 117 million parameters, are trained on the BooksCorpus dataset (800 million words).

# GPT-2: WebText Dataset

## WebText Dataset:

- Curated using Reddit as a **filter** to include diverse, high-quality content.
- Comprised of approximately 8 million documents and 40GB of text.

## Training:

- Trained as a **standard language model** on WebText.
- Objective: Maximize the likelihood of predicting the next token in a sequence.

## Zero-Shot Inference:

- Inputs during inference are formatted as **prompts** that specify tasks.
- Generates responses based on pre-trained language understanding, without fine-tuning.
- Example: **Input:** "Q: What is the capital of France? A:" and **Response:** "Paris"

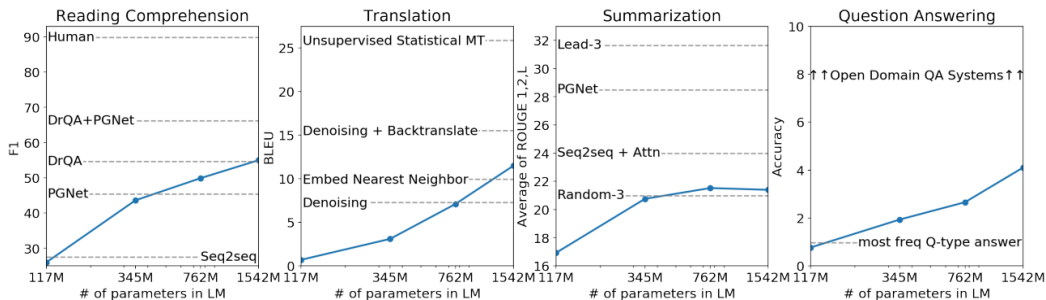
## Model Details:

- 12-layer Transformer decoder with 1.5 billion parameters.
- Hidden size of 1600, with 25 attention heads.

# GPT-2: Zero-Shot Inference

Parameters	Layers	$d_{model}$
117M	12	768
345M	24	1024
762M	36	1280
1542M	48	1600

**Model configurations**



**Zero-shot performance of WebText GPT-2 on many NLP tasks.**

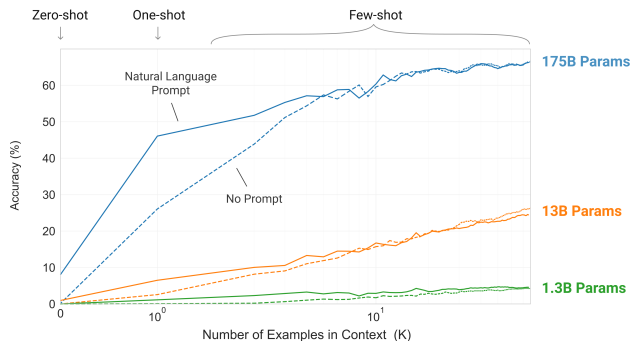
# GPT-3: Limitation of Pretraining and Finite-Tune

## Limitation of Pretrained and Finite-Tune in NLP:

- **Large Data Limitation:** Language models need extensive **labeled** datasets for new tasks, limiting their practical use.
- **Overfitting and Generalization Issues:** Powerful models risk overfitting, especially when fine-tuned on small data with limited diversity, leading to poor real-world generalization.
- **Human Adaptability:** Humans learn language tasks with minimal examples and switch tasks seamlessly; achieving this in NLP systems would improve their versatility.

## Solution: In-context learning and Transformer

- **In-Context Learning:** the model is conditioned on a natural language instruction and a few **demonstrations** of the task and is then expected to complete further instances of the task simply by predicting what comes next.





# GPT-3: In-Context Learning

The three settings we explore for in-context learning

## Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 cheese => ..... ← prompt
```

## One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 sea otter => loutre de mer ← example
3 cheese => ..... ← prompt
```

## Few-shot

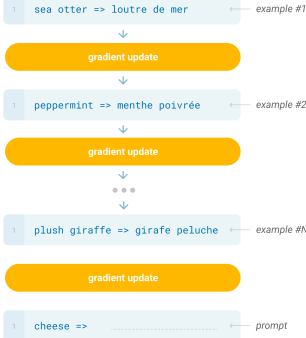
In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 sea otter => loutre de mer ← examples
3 peppermint => menthe poivrée ←
4 plush girafe => girafe peluche ←
5 cheese => ..... ← prompt
```

Traditional fine-tuning (not used for GPT-3)

## Fine-tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.



## GPT-3: Training

Model Name	$n_{\text{params}}$	$n_{\text{layers}}$	$d_{\text{model}}$	$n_{\text{heads}}$	$d_{\text{head}}$	Batch Size	Learning Rate
GPT-3 Small	125M	12	768	12	64	0.5M	$6.0 \times 10^{-4}$
GPT-3 Medium	350M	24	1024	16	64	0.5M	$3.0 \times 10^{-4}$
GPT-3 Large	760M	24	1536	16	96	0.5M	$2.5 \times 10^{-4}$
GPT-3 XL	1.3B	24	2048	24	128	1M	$2.0 \times 10^{-4}$
GPT-3 2.7B	2.7B	32	2560	32	80	1M	$1.6 \times 10^{-4}$
GPT-3 6.7B	6.7B	32	4096	32	128	2M	$1.2 \times 10^{-4}$
GPT-3 13B	13.0B	40	5140	40	128	2M	$1.0 \times 10^{-4}$
GPT-3 175B	175.0B	96	12288	96	128	3.2M	$0.6 \times 10^{-4}$

Dataset	Quantity (tokens)	Weight in training mix	Epochs elapsed when training for 300B tokens
Common Crawl (filtered)	410 billion	60%	0.44
WebText2	19 billion	22%	2.9
Books1	12 billion	8%	1.9
Books2	55 billion	8%	0.43
Wikipedia	3 billion	3%	3.4

# GPT-3: Neural Scaling Laws

