

Sequence-to-Sequence Models

Tianxiang (Adam) Gao

February 27, 2025

Outline

- 1 Sequence-to-Sequence Models
- 2 Attention Mechanism
- 3 Transformer
- 4 Large Language Models: BERT

Recap: Recurrent Neural Networks

- **Audio Waveform:** A 1D array represents the amplitude of the sound over time, e.g., 16kHz
- **One-Hot Encoding:** Each word in a **vocabulary** is a **binary** one-hot vector.
- **Challenges in Text Data:** Curse of dimensionality and long-run dependencies.
- **Language Models:** Assigns probabilities to a given sequence of words
- **Neural Language Model:** Model the probability distribution of the next word given the history:

$$\mathbb{P}(\mathbf{x}_{t+1} \mid \mathbf{x}_1, \dots, \mathbf{x}_t) = f_{\theta}(\mathbf{x}_1, \dots, \mathbf{x}_t).$$

RNNs: Encode the history into a hidden state \mathbf{h}_t updated by combining with the current word \mathbf{x}_t :

$$\mathbf{h}_t = \tanh(\mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{W}_x \mathbf{x}_t + \mathbf{b}_h), \quad \text{and} \quad \hat{\mathbf{y}}_t = \text{softmax}(\mathbf{W}_y \mathbf{h}_t + \mathbf{b}_y)$$

Training RNNs: backpropagation through time

- Forward (simplified): $\mathbf{h}_t = \phi(\mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{W}_x \mathbf{x}_t)$
- Backward (simplified): $d\mathbf{h}_t = \mathbf{W}_h^{\top} (\phi'_{t+1} \odot d\mathbf{h}_{t+1}) + \mathbf{W}_y^{\top} (\sigma'_t \odot d\mathbf{y}_t)$
- **Generation:** Sample the next word from the predicted probability distribution produced by RNNs.
- **RNN Types:** One-to-many, many-to-one, or many-to-many structures for different tasks.
- **Vanishing/Exploding Gradients:** $\mathbf{h}_t = \mathcal{O}(a^t)$ and $d\mathbf{h}_t = \mathcal{O}(b^{T-t})$.

Recap: Recurrent Neural Networks

- **Gated Recurrent Unit (GRU):** Gates helps maintain long-term dependencies:

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_h(\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{W}_x \mathbf{x}_t), \quad \text{and} \quad \mathbf{h}_t = \mathbf{z}_t \mathbf{h}_{t-1} + (1 - \mathbf{z}_t) \odot \tilde{\mathbf{h}}_t$$

Long Short-Term Memory (LSTM): Use a cell state \mathbf{c}_t to maintain long-term dependencies.

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{c}}_t, \quad \text{and} \quad \mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)$$

Bidirectional RNNs: The concatenated hidden state: $\mathbf{h}_t = [\vec{\mathbf{h}}_t, \overleftarrow{\mathbf{h}}_t]$

- **Deeper RNNs:** Each layer ℓ computes its hidden state using the hidden state from the layer $\ell - 1$:

$$\mathbf{h}_t^{(\ell)} = \tanh(\mathbf{W}_h^{(\ell)} \mathbf{h}_{t-1}^{(\ell)} + \mathbf{W}_x^{(\ell)} \mathbf{x}_t^{(\ell-1)})$$

Drawbacks of One-Hot Representation: orthogonality and high dimensionality

- **Word Embedding:** Words are represented as **dense vectors** in a **lower-dimensional space**.

$$\mathbf{e} = \mathbf{E}\mathbf{x}$$

Continuous Bag of Words (CBOW): Predicts the target word given the context

- **Skip-Gram:** Predicts the context words given a target word.
- **Negative Sampling:** Reformulates the context-target predictions as a **binary** classification:

$$\mathcal{L}(\mathbf{E}) = - \sum_{(\mathbf{e}_c, \mathbf{e}_t)} \log \sigma(\mathbf{e}_t^\top \mathbf{e}_c) - \sum_{(\mathbf{e}_c, \tilde{\mathbf{e}}_t)} \log \sigma(-\mathbf{e}_t^\top \tilde{\mathbf{e}}_c),$$

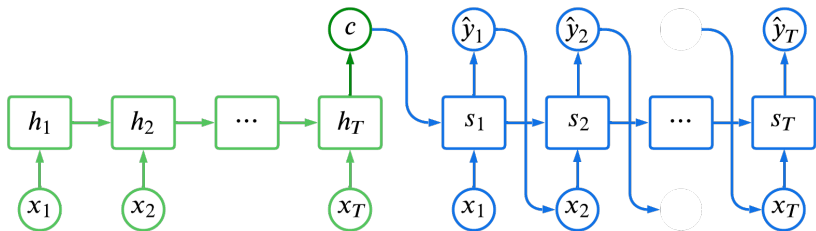
where $\tilde{\mathbf{e}}_t$ is **negative** target samples outside the context window.

Outline

- 1 Sequence-to-Sequence Models
- 2 Attention Mechanism
- 3 Transformer
- 4 Large Language Models: BERT

Seq2Seq: Sequence-to-Sequence Models

- **Define:** Sequence-to-Sequence (Seq2Seq) models are designed to handle tasks where both input and output are sequences of **variable** length, e.g., machine translation or summarization.
- **Example:** “Jane visite l’Afrique en septembre.” \implies “Jane visits Africa in September.”
- **Encoder-Decoder Architecture**
 - **Encoder:** Processes and compresses the input sequence into a fixed-length **context vector**.
 - **Decoder:** Uses the context vector to generate the output sequence sequentially



Here c is the context vector summarizing the *entire* input sequence.

Conditional Language Model

- **Language Model:** Assigns probabilities to a sequence of words $\{\mathbf{x}_t\} = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$:

$$\mathbb{P}(\{\mathbf{x}_t\}) = \mathbb{P}(\mathbf{x}_1, \dots, \mathbf{x}_T) = \mathbb{P}(\mathbf{x}_1) \cdot \prod_{t=2}^{T-1} \mathbb{P}(\mathbf{x}_{t+1} \mid \mathbf{x}_1, \dots, \mathbf{x}_t)$$

where each conditional probability is modeled as:

$$\mathbb{P}(\mathbf{x}_{t+1} \mid \mathbf{x}_1, \dots, \mathbf{x}_t) = f_{\theta}(\mathbf{x}_1, \dots, \mathbf{x}_t)$$

Conditional Language Model: Assigns probabilities to a target sequence $\{\mathbf{y}_t\}$ *given* an input sequence $\{\mathbf{x}_t\}$:

$$\mathbb{P}(\{\mathbf{y}_t\} \mid \{\mathbf{x}_t\}) = \prod_{t=1}^{T'-1} \mathbb{P}(\mathbf{y}_{t+1} \mid \mathbf{y}_1, \dots, \mathbf{y}_t, \mathbf{x}_1, \dots, \mathbf{x}_T)$$

where the conditional probability for each word in the target sequence is:

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

with $(\mathbf{x}_1, \dots, \mathbf{x}_T) \mapsto \mathbf{c}$ and $(\mathbf{y}_1, \dots, \mathbf{y}_{t-1}) \mapsto \mathbf{s}_{t-1}$.

Beam Search

Greedy Search: Selects the most probable word at each step, which may lead to suboptimal sequences.

Beam Search: Tracks **multiple** high-probability sequences simultaneously to improve overall accuracy.

- **Initialization:** Start with the *seed* token and choose the top k words based on the probability distribution.
- **Expansion:** Predict the next word for each candidate, generating new sequences.
- **Pruning:** Keep the top k sequences with the highest log-probability scores, discarding the rest.

$$\begin{aligned}\log \mathbb{P}(\mathbf{y}_2, \mathbf{y}_1 \mid \mathbf{c}) &= \log \left[\mathbb{P}(\mathbf{y}_2 \mid \mathbf{c}, \mathbf{y}_1) \cdot \mathbb{P}(\mathbf{y}_1 \mid \mathbf{c}) \right] \\ &= \log \mathbb{P}(\mathbf{y}_2 \mid \mathbf{c}, \mathbf{y}_1) + \log \mathbb{P}(\mathbf{y}_1 \mid \mathbf{c}).\end{aligned}$$

Repeat: Continue expanding and pruning until terminate

Remark

- **Beam Width (k):** Requires k identical decoders to update candidate sequences simultaneously.
- **Advantages:** Balances between accuracy and computation; larger k increases accuracy but demands more computational resources.

Numerical Stability

Log-Probabilities for Stability:

- Since $\mathbb{P}(\cdot) \in [0, 1]$, the product of probabilities can approach zero, causing numerical instability.
- To address this, log-probabilities are used:

$$\begin{aligned}\mathbf{y}^* &= \operatorname{argmax}_{\mathbf{y}} \mathbb{P}(\mathbf{y}_1, \dots, \mathbf{y}_{T'} \mid \mathbf{c}) \\ &= \operatorname{argmax}_{\mathbf{y}} \frac{1}{T'} \sum_{t=1}^{T'} \log \mathbb{P}(\mathbf{y}_t \mid \mathbf{c}, \mathbf{y}_1, \dots, \mathbf{y}_{t-1})\end{aligned}$$

This transformation helps prevent underflow and enables stable computation of probabilities.

Error Analysis in Beam Search

Let \mathbf{y}^* be the optimal sequence and $\hat{\mathbf{y}}$ the model's predicted sequence.

$$\mathbf{y}^* = \underset{\mathbf{y}}{\operatorname{argmax}} \mathbb{P}(\mathbf{y} \mid \mathbf{c}) = \mathbb{P}(\mathbf{y}_1, \dots, \mathbf{y}_{T'} \mid \mathbf{c})$$

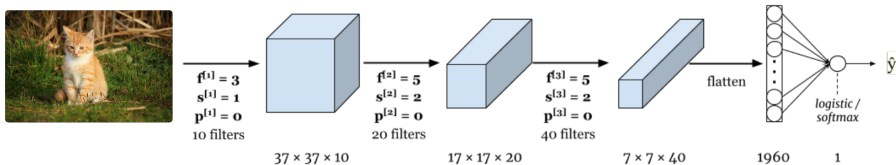
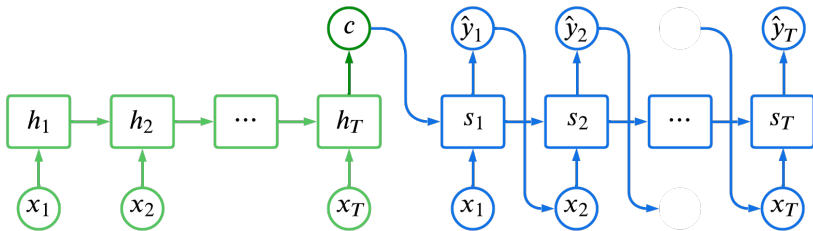
If $\mathbb{P}(\mathbf{y}^* \mid \mathbf{c}) > \mathbb{P}(\hat{\mathbf{y}} \mid \mathbf{c})$, then

- Beam search pick $\hat{\mathbf{y}}$ but \mathbf{y}^* is better, achieving higher $p(\mathbf{y} \mid \mathbf{c})$
- Increasing beam width can improve accuracy by exploring more potential sequences

If $\mathbb{P}(\mathbf{y}^* \mid \mathbf{c}) \leq \mathbb{P}(\hat{\mathbf{y}} \mid \mathbf{c})$, then

- \mathbf{y}^* is a better translation than $\hat{\mathbf{y}}$, but RNN predicts $\mathbb{P}(\hat{\mathbf{y}} \mid \mathbf{c}) \geq \mathbb{P}(\mathbf{y}^* \mid \mathbf{c})$
- RNN module is at fault
- Increasing beam width won't help, as errors stem from model limitations.

Image Captioning: Using CNN as the Encoder



"A small orange kitten sits attentively on green grass, surrounded by natural, dried foliage in the background, giving a calm and serene outdoor setting."

BLEU Score: Bilingual Evaluation Understudy

BLEU Score: Bilingual Evaluation Understudy (BLEU) evaluates a machine-generated **candidate translation** \hat{y} by comparing it to a *list* of reference translations $\{y_1, \dots, y_M\}$.

- **Candidate Translation:** "the cat the cat on the mat"
- **Reference Translation 1:** "the cat is on the mat"
- **Reference Translation 2:** "there is a cat on the mat"

Precision: Measures how many n-grams in the candidate match the reference translations.

- **Candidate Bigrams:** {"the cat", "cat the", "the cat", "cat on", "on the", "the mat"}
- **Reference 1 Bigrams:** {"the cat", "cat is", "is on", "on the", "the mat"}
- **Reference 2 Bigrams:** {"there is", "is a", "a cat", "cat on", "on the", "the mat"}

$$\text{Precision} = \frac{\text{Number of matching n-grams}}{\text{Total n-grams in candidate}} = \frac{5}{6} \implies \text{Inflate the score!}$$

Modified Precision

Modified Precision: Limits the count of an n-gram to the maximum it appears in any reference.

- **Candidate Bigrams:** {"the cat", "cat the", "the cat", "cat on", "on the", "the mat"}
- **Reference 1 Bigrams:** {"the cat", "cat is", "is on", "on the", "the mat"}
- **Reference 2 Bigrams:** {"there is", "is a", "a cat", "cat on", "on the", "the mat"}

Bigrams	Count	Matching	Clipped
the cat	2	2	1
cat the	1	0	0
cat on	1	2	1
on the	1	3	1
the mat	1	2	1

The corresponding modified precision is given by:

$$\text{Modified Precision} = \frac{\text{Clipped number of matching n-grams}}{\text{Total n-grams in candidate}} = \frac{4}{6}$$

BLEU Score Formula

- For an n -gram, the **modified precision** p_n is defined as:

$$p_n = \frac{\text{Clipped number of matching } n\text{-grams}}{\text{Total } n\text{-grams in candidate}}$$

- The **BLEU score** is computed using the average of modified precisions, combined with a **brevity penalty (BP)** to penalize overly short translations:

$$\text{BLEU} = \text{BP} \cdot \exp\left(\frac{1}{N} \sum_{n=1}^N \log p_n\right)$$

where N is the maximum n -gram length considered (typically $N = 4$).

- The **Brevity Penalty (BP)** prevents high scores for short translations:

$$\text{BP} = \begin{cases} 1, & \text{if } c \geq r \\ \exp\left(1 - \frac{r}{c}\right), & \text{if } c < r \end{cases}$$

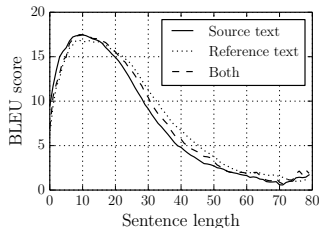
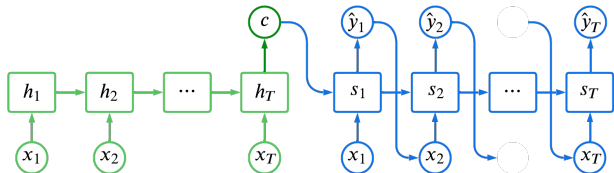
where:

- c is the length of the candidate translation.
- r is the length of the reference translation.

Outline

- 1 Sequence-to-Sequence Models
- 2 Attention Mechanism**
- 3 Transformer
- 4 Large Language Models: BERT

Limitations of RNN Encoder-Decoder Framework



- **Encoder:** Process and compress the input sequence $\{x_1, \dots, x_T\}$ into a context vector c .

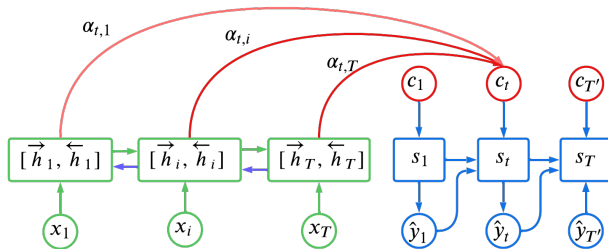
$$h_t = f_\theta(h_{t-1}, x_t), \quad c = h_T.$$

- **Decoder:** Uses the context vector c to generate the output sequence sequentially

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

Note: Encoding all information into a **single** vector c may cause information loss for longer sequences.

Distinct Context Vector in Attention Mechanism



- **Distinct Context Vector for Each Target Word:** Each target word y_t has a unique context vector c_t , allowing the model to focus on relevant input parts.

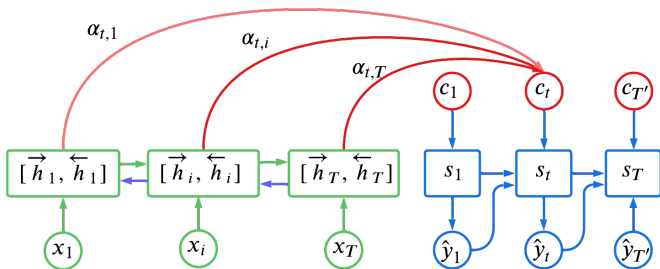
$$\mathbb{P}(y_t | x, y_1, \dots, y_{t-1}) = \mathbb{P}(y_t | s_t), \quad \text{where } s_t = g_\phi(s_{t-1}, y_{t-1}, c_t)$$

Context Vector Computation: The context vector c_t is computed as a weighted sum of encoder hidden states h_i , tailored to the current decoding step.

$$c_t = \sum_{i=1}^T \alpha_{t,i} h_i$$

where **attention weights** $\alpha_{t,i}$ indicate the relevance of each hidden state h_i for generating y_t .

Attention Mechanism in Seq2Seq



- **Attention Weights:** Computed from an **alignment score** $e_{t,i}$ between the decoder's previous hidden state s_{t-1} and each encoder hidden state h_i .

$$\alpha_{t,i} = \frac{e_{t,i}}{\sum_j e_{t,j}} \quad \text{where} \quad e_{t,i} = a(s_{t-1}, h_i) = \mathbf{v}^\top \tanh(\mathbf{W}s_{t-1} + \mathbf{U}h_i)$$

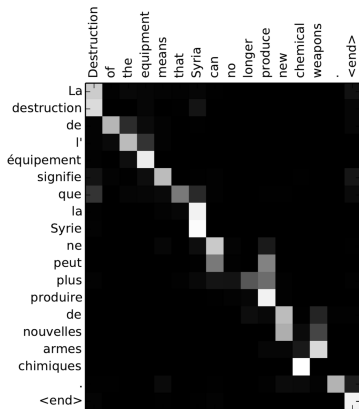
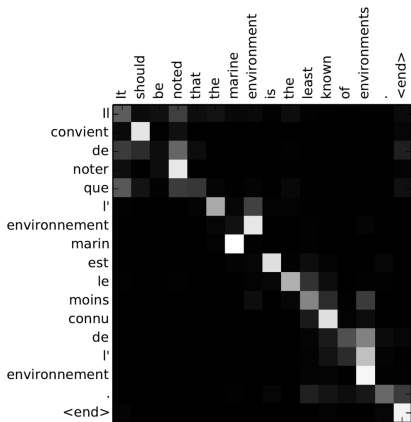
is an **alignment model** implemented as an MLP, which is trained jointly with Seq2Seq.

- **Bidirectional Encoder:** The encoder uses a bidirectional RNN to capture both past and future context, enhancing comprehension of each input word's meaning.

Sample Alignments

- **Attention Weights:** Computed from an **alignment score** $e_{t,i}$ between the decoder's previous hidden state s_{t-1} and each encoder hidden state h_i .

$$\alpha_{t,i} = \frac{e_{t,i}}{\sum_j e_{t,j}} \quad \text{where} \quad e_{t,i} = a(s_{t-1}, h_i)$$



Summary: Seq2Seq Models

- **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

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

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

$$\mathbf{c}_t = \sum_i \alpha_{t,i} \mathbf{h}_i,$$

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

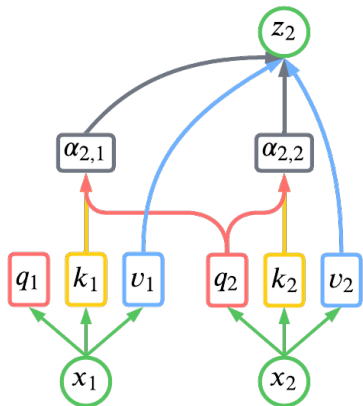
- **Alignment Scores:** Alignment scores $e_{t,i} = a(\mathbf{s}_{t-1}, \mathbf{h}_i)$ indicate relevance between encoder hidden states \mathbf{h}_i and decoder states \mathbf{s}_{t-1} .

Outline

- 1 Sequence-to-Sequence Models
- 2 Attention Mechanism
- 3 Transformer**
- 4 Large Language Models: BERT

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

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

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

$$\mathbf{Q} = [q_1 \quad \cdots \quad q_T]^\top, \quad \mathbf{K} = [k_1 \quad \cdots \quad k_T]^\top, \quad \mathbf{V} = [v_1 \quad \cdots \quad v_T]^\top$$

The attention weights are computed as:

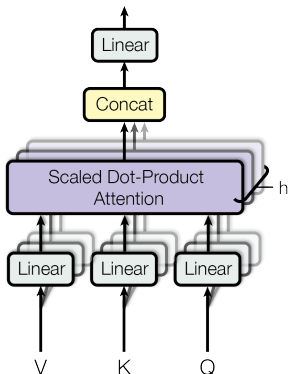
$$[\alpha_1 \quad \cdots \quad \alpha_T] = \text{softmax} \left(\frac{\mathbf{KQ}^\top}{\sqrt{d_k}} \right)$$

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:

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

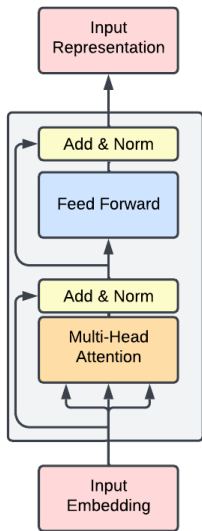
where W_o is the output projection and H is the number of heads.

- **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 $\mathcal{O}(T^2d)$.

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:

$$z = \mathbf{W}x, \quad \mathbf{h} = \tanh(z)$$

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

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

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

where \mathbf{w}_i is the i th row of \mathbf{W} .

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

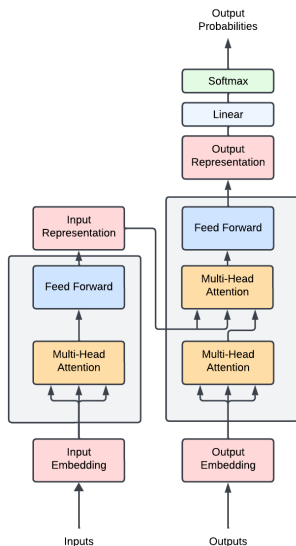
$$z_{\text{norm}} = \frac{z - \mu}{\sigma}, \quad \tilde{z} = \alpha \odot z_{\text{norm}} + \beta$$

where α and β are **trainable**.

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

$$\text{FFN}(x_t) = \text{ReLU}(x_t \mathbf{W}_1 + \mathbf{b}_1) \mathbf{W}_2 + \mathbf{b}_2$$

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.

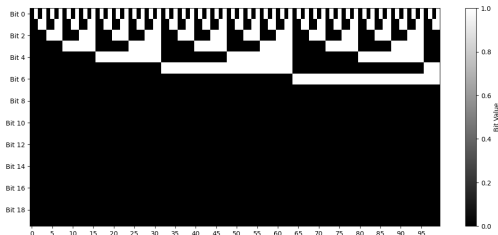
Decoder:

- **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.

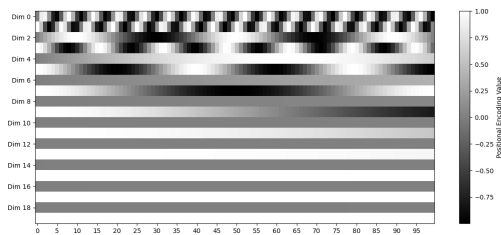
Positional Encoding

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

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



Binary Representation

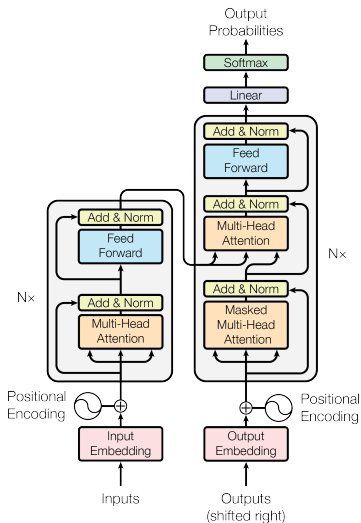


Position Encoding

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. $M_k PE_{\text{pos}+k} = PE_{\text{pos}}$.

Transformer



Training Process: Teacher Forcing

- During training, the **ground truth (actual) previous token** 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

- 1 Sequence-to-Sequence Models
- 2 Attention Mechanism
- 3 Transformer
- 4 Large Language Models: BERT**

BERT

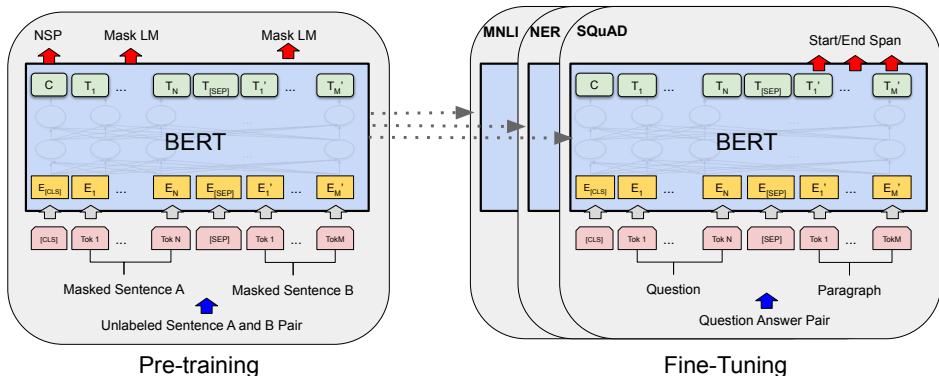


Google BERT

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 _{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