

NLP: lab2

序号	学号	专业班级	姓名	性别
/	3220102157	图灵2201	张祎迪	女

1.Project Introduction

Development Environment:

- ModelArts Ascend Notebook : mindspore1.7.0-cann5.1.0-py3.7-euler2.8.3

SeqtoSeq with Transformers

Sequence-to-Sequence (Seq2Seq) Model

- Seq2Seq Model:** A type of neural network architecture used for sequence prediction tasks, such as machine translation, text summarization, and speech recognition.
- Encoder-Decoder Architecture:** Consists of an encoder network that processes the input sequence and a decoder network that generates the output sequence.
- Attention Mechanism:** Allows the model to focus on different parts of the input sequence during the decoding process.

This project aims to implement a Seq2Seq model using the Transformer architecture for machine translation tasks. The Transformer model is a state-of-the-art architecture that has been widely adopted for sequence-to-sequence tasks due to its parallelism and scalability.

Steps:

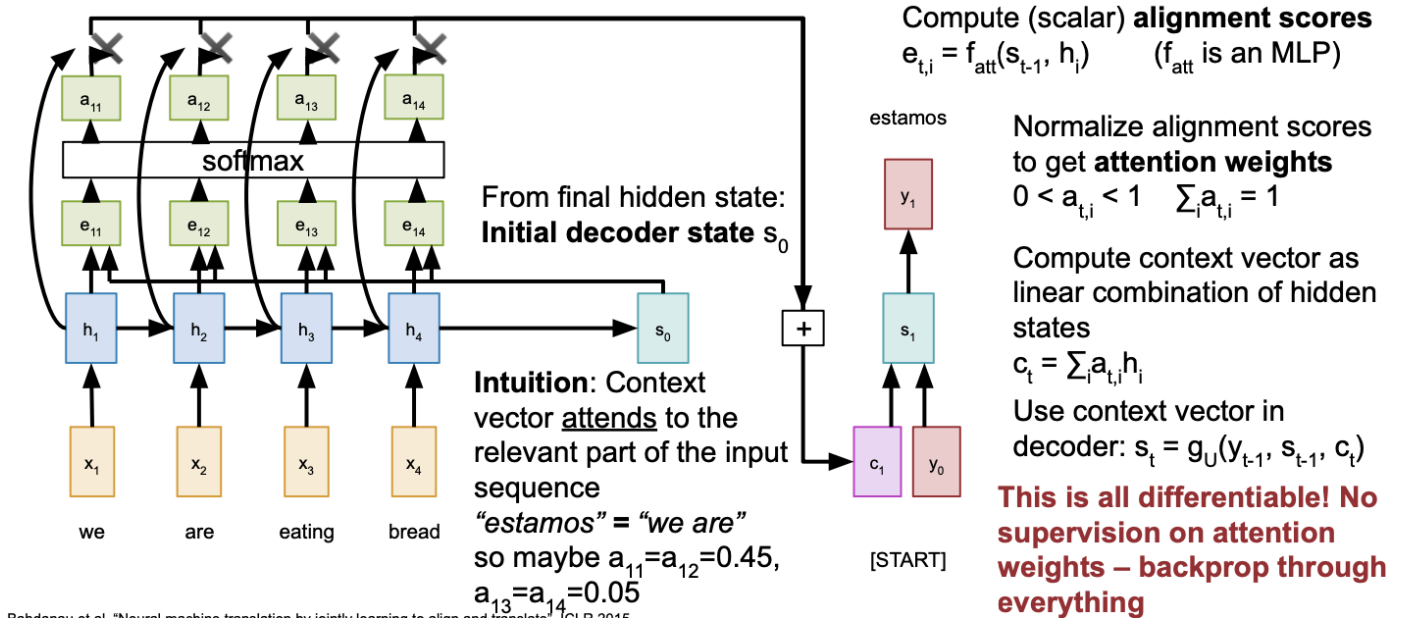
- Data Preprocessing:** Prepare the training and validation datasets for the machine translation task.
- Model Architecture:** Implement the Transformer model architecture for sequence-to-sequence tasks.
- Training:** Train the model on the training dataset and evaluate its performance on the validation dataset.
- Inference:** Use the trained model to generate translations for new input sequences.

2.Technical Details

Theoretically Elaboration

Attention Mechanism

Basic Idea: The attention mechanism allows the model to focus on different parts of the input sequence during the decoding process. It assigns weights to the input elements based on their relevance to the current output element, enabling the model to capture long-range dependencies and improve the quality of the generated sequences.



From the above figure, we can see that each s_i is computed based on the weighted sum of the input elements h_1, h_2, \dots, h_n , where the weights are determined by the attention scores a_{ij} .

Improved Attention Mechanisms

Self-Attention Layer

Inputs:

Input vectors: X (Shape: $N_x \times D_x$)

Key matrix: W_K (Shape: $D_x \times D_Q$)

Value matrix: W_V (Shape: $D_x \times D_V$)

Query matrix: W_Q (Shape: $D_x \times D_Q$)

Computation:

Query vectors: $Q = XW_Q$

Key vectors: $K = XW_K$ (Shape: $N_x \times D_Q$)

Value Vectors: $V = XW_V$ (Shape: $N_x \times D_V$)

Similarities: $E = QK^T / \sqrt{D_Q}$ (Shape: $N_x \times N_x$) $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$

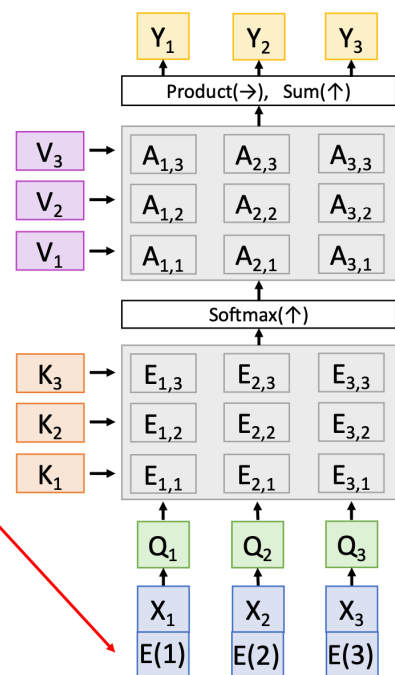
Attention weights: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_x \times N_x$)

Output vectors: $Y = AV$ (Shape: $N_x \times D_V$) $Y_i = \sum_j A_{i,j} V_j$

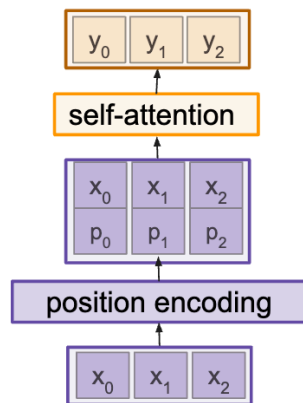
Self attention doesn't “know” the order of the vectors it is processing!

In order to make processing position-aware, concatenate or add **positional encoding** to the input

E can be learned lookup table, or fixed function



- **Self-Attention:** A mechanism that allows the model to focus on different parts of the input sequence during the decoding process.
- **Scaled Dot-Product Attention:** A type of self-attention mechanism that scales the dot-product attention scores by the square root of the dimensionality of the key vectors.



Concatenate special positional encoding p_j to each input vector x_j

We use a function $pos: N \rightarrow R^d$ to process the position j of the vector into a d -dimensional vector

So, $p_j = pos(j)$

Options for $pos(\cdot)$

1. Learn a lookup table:
 - Learn parameters to use for $pos(t)$ for $t \in [0, T]$
 - Lookup table contains $T \times d$ parameters.
2. Design a fixed function with the desiderata

Intuition:

$$p(t) = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \sin(\omega_2 \cdot t) \\ \cos(\omega_2 \cdot t) \\ \vdots \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_d$$

0 :	0	0	0	0
1 :	0	0	0	1
2 :	0	0	1	0
3 :	0	0	1	1
4 :	0	1	0	0
5 :	0	1	0	1
6 :	0	1	1	0
7 :	0	1	1	1
8 :	1	0	0	0
9 :	1	0	0	1
10 :	1	0	1	0
11 :	1	0	1	1
12 :	1	1	0	0
13 :	1	1	0	1
14 :	1	1	1	0
15 :	1	1	1	1

where $\omega_k = \frac{1}{10000^{2k/d}}$

[image](#)
Vaswani et al, "Attention is all you need", NeurIP

- **Positional Encoding:** A technique used to inject positional information into the input embeddings to capture the sequential order of the input sequence.

Multihead Self-Attention

Inputs:

Input vectors: X (Shape: $N_X \times D_X$)

Key matrix: W_K (Shape: $D_X \times D_Q$)

Value matrix: W_V (Shape: $D_X \times D_V$)

Query matrix: W_Q (Shape: $D_X \times D_Q$)

Use H independent "Attention Heads" in parallel

Computation:

Query vectors: $Q = XW_Q$

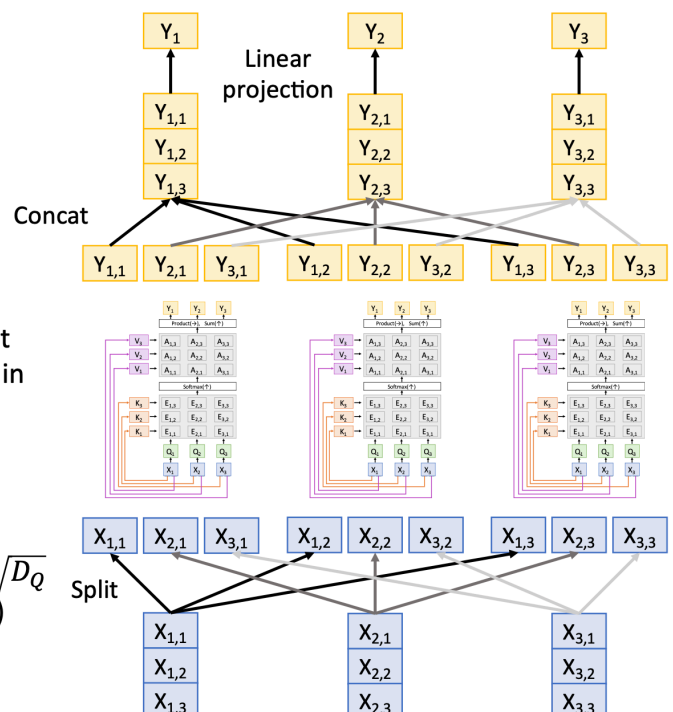
Key vectors: $K = XW_K$ (Shape: $N_X \times D_Q$)

Value Vectors: $V = XW_V$ (Shape: $N_X \times D_V$)

Similarities: $E = QK^T / \sqrt{D_Q}$ (Shape: $N_X \times N_X$) $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$

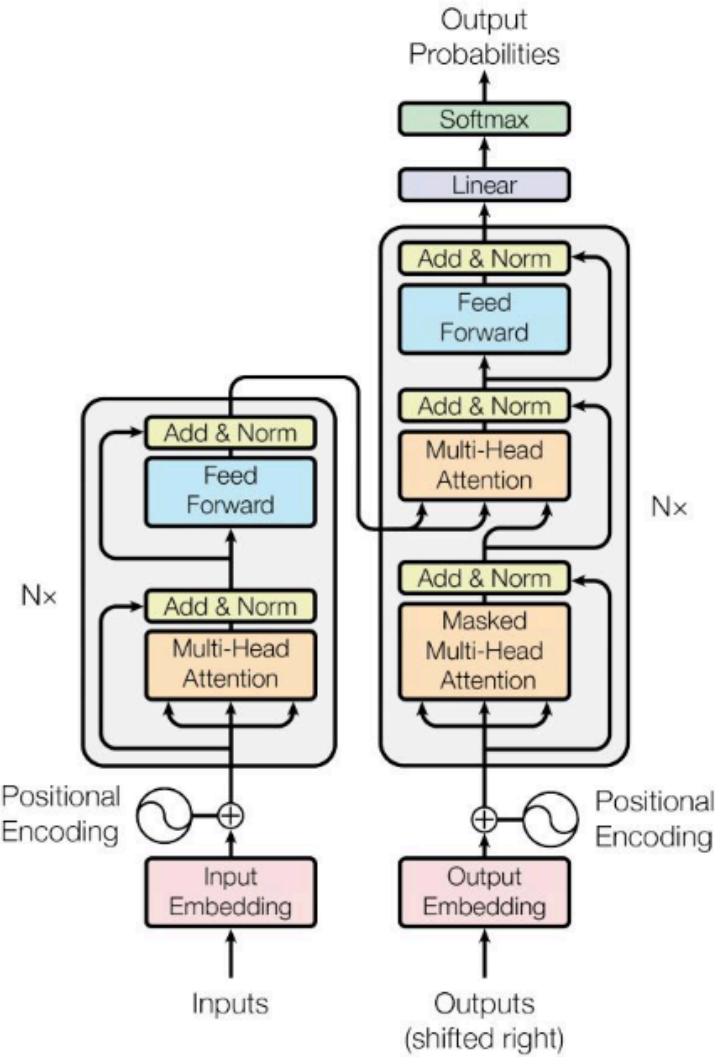
Attention weights: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_X \times N_X$)

Output vectors: $Y = AV$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} V_j$



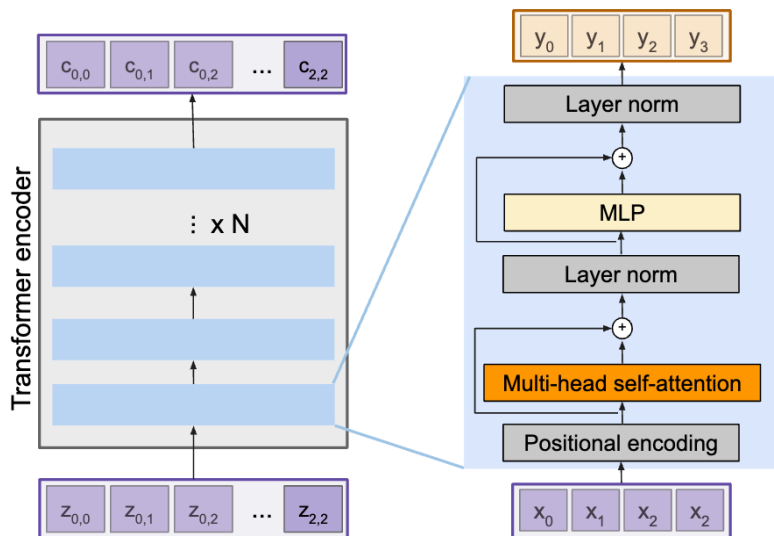
- **Multi-Head Attention:** A mechanism that allows the model to jointly attend to information from different representation subspaces.

The Transformer Architecture



- **Encoder:** Processes the input sequence and generates a sequence of hidden representations.

The Transformer encoder block



Transformer Encoder Block:

Inputs: Set of vectors \mathbf{x}

Outputs: Set of vectors \mathbf{y}

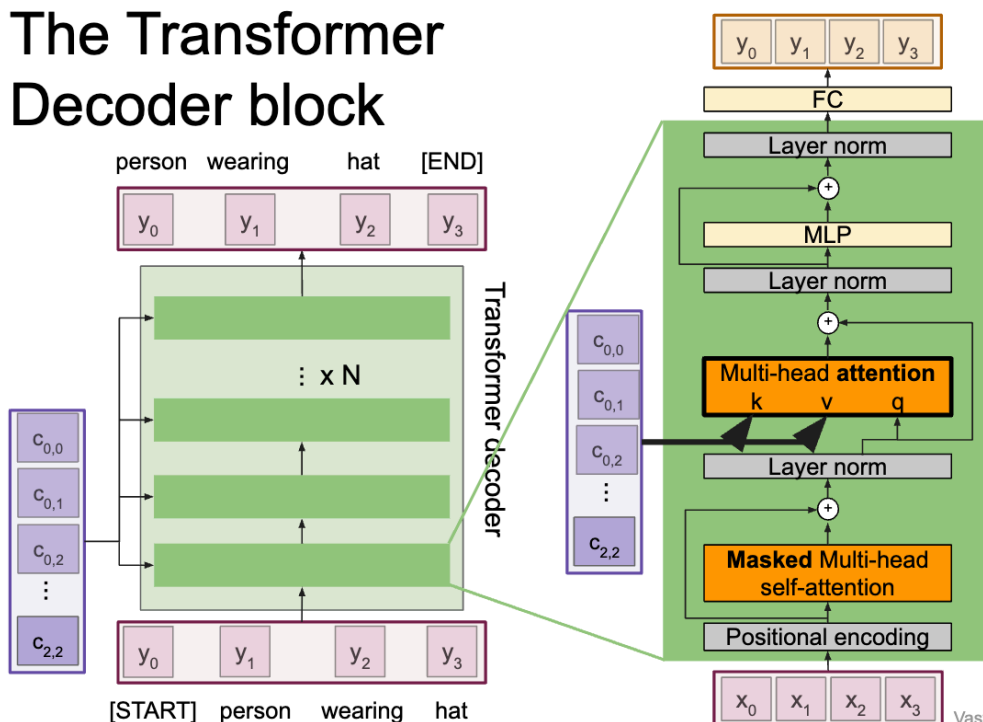
Self-attention is the only interaction between vectors.

Layer norm and MLP operate independently per vector.

Highly scalable, highly parallelizable, but high memory usage.

- **Decoder:** Generates the output sequence based on the encoder's hidden representations and the previous output elements.

The Transformer Decoder block



Multi-head attention block attends over the transformer encoder outputs.

For image captioning, this is how we inject image features into the decoder.

Vaswani et al, "Attention is all you need", NeurIPS 2017

Algorithm Implementation

I did not implement the algorithm myself in the project, but I used the MindSpore library to build the Transformer model for machine translation tasks. The MindSpore library provides a high-level API for building and training deep learning models, making it easy to implement complex architectures such as the Transformer model.

1. Data Preprocessing

```
1 sample_num = 23607
2 eval_idx = np.random.choice(sample_num, int(sample_num*0.2), replace=False)
3 data_prepare(data_cfg, eval_idx)
```

Choose 20% of the data as training set

The `data_prepare` function preprocesses input data, tokenizes it, filters instances based on length constraints, and writes the preprocessed data to separate files for training and evaluation in a format suitable for training a transformer-based model.

2. Training

- Main part of the parameters (in the `train_cfg` dictionary):

```
'''
two kinds of transformer model version
'''
if train_cfg.transformer_network == 'base':
    transformer_net_cfg = TransformerConfig(
        batch_size=train_cfg.batch_size,
        seq_length=train_cfg.seq_length,
        vocab_size=train_cfg.vocab_size,
        hidden_size=512,
        num_hidden_layers=6,
        num_attention_heads=8,
        intermediate_size=2048,
        hidden_act="relu",
        hidden_dropout_prob=0.2,
        attention_probs_dropout_prob=0.2,
        max_position_embeddings=train_cfg.max_position_embeddings,
        initializer_range=0.02,
        label_smoothing=0.1,
        input_mask_from_dataset=True,
        dtype=mstype.float32,
        compute_type=mstype.float16)
elif train_cfg.transformer_network == 'large':
    transformer_net_cfg = TransformerConfig(
        batch_size=train_cfg.batch_size,
        seq_length=train_cfg.seq_length,
        vocab_size=train_cfg.vocab_size,
        hidden_size=1024,
        num_hidden_layers=6,
        num_attention_heads=16,
        intermediate_size=4096,
        hidden_act="relu",
        hidden_dropout_prob=0.2,
        attention_probs_dropout_prob=0.2,
        max_position_embeddings=train_cfg.max_position_embeddings,
        initializer_range=0.02,
        label_smoothing=0.1,
        input_mask_from_dataset=True,
        dtype=mstype.float32,
        compute_type=mstype.float16)
else:
    raise Exception("The src/train_config of transformer_network must base or large. Change the src/train_config file and try ag
```

Specifications

`batch_size`: The number of instances in each batch.

`seq_length`: The maximum sequence length for input and output sequences.

`vocab_size`: The size of the vocabulary used for tokenization.

`hidden_size`: The size of the hidden layers in the Transformer model.

`num_hidden_layers`: The number of hidden layers in the Transformer model.

`num_attention_heads`: The number of attention heads in the Transformer model.

`hidden_act`: The activation function used in the hidden layers.

`hidden_dropout_prob`: The dropout probability for the hidden layers.

`attention_probs_dropout_prob`: The dropout probability for the attention probabilities.

`max_position_embeddings`: The maximum number of positions for positional embeddings.

`initializer_range`: The range for random weight initialization.

`label_smoothing`: The label smoothing factor for the loss function.

`input_mask_from_dataset`: Whether to use input masks from the dataset.

```
1 | train(train_cfg)
```

The `train()` function trains the model. It first loads the dataset and initializes the Transformer network, loss function, learning rate, and optimization method. If pre-trained model parameters exist, it loads the parameters from the specified `cfg.checkpoint_path` path into the network. If model saving is enabled in the configuration, it also adds a `ModelCheckpoint` callback function to save the model. Additionally, it adjusts the optimizer based on whether `enable_lossscale` is enabled. Finally, it wraps everything into a model using the `Model` class, sets it to training mode, and calls the `train()` method to start training.

3. Evaluation

```
1 | evaluate(eval_cfg)
```

- Specific parameters in the `eval_cfg` dictionary
- More details can be found in the `mt_transformer_mindspore.ipynb` file.

3.Experiment Results

- All test results can be found in the `mt_transformer_mindspore.ipynb` file.

```
...
result: 他有可能再迟到了。
source: He is lying through his teeth . 他明显在撒谎。
result: 他穿过衣服的牙齿穿衣服在牙齿。
source: He painted a picture of roses . 他画了一幅玫瑰的画。
result: 他画了一幅画的照片。
source: He put on an air of innocence . 他摆出一副无辜的样子。
result: 他在空空空空空空空气上发现了空气。
source: He studies history at college . 他在大学修读历史。
result: 他在大学校学的历史学习历史上学习历史书。
source: He threw a rock into the pond . 他扔一块石头到池塘里。
result: 他扔了一块岩石头后扔了。
source: He was busy with his homework . 他忙于做功课。
result: 他忙于做他的家庭作业后做作业。
source: He was hurt in a car accident . 他在一次车祸中受伤了。
result: 他在一场车祸中发生了意外的意外。
source: He ' ll be back by five o ' clock . 他五点左右会回来。
result: 他五点之前会五点回来。
source: He ' s a friend of my brother ' s . 他是我哥哥的朋友。
result: 他是我弟弟弟弟弟弟的一个朋友。
source: His wife is one of my friends . 他的妻子是我的一个朋友。
result: 他的妻子是我的一个朋友都是我的朋友。
source: How about going to the movies ? 我们去电影院怎么样？
result: 去看电影怎么样？
source: How are you going to get home ? 你打算怎么回家？
result: 你要怎么去家里？
source: How much is this handkerchief ? 请问这个手帕多少钱？
result: 今天的手帕怎么样？
source: I am interested in this story . 我对这个故事感兴趣。
result: 我对这个故事对这个故事感兴趣。
```

- Some typical results are shown below:

Lack in learning slang

```
1 source: He is lying through his teeth . 他明显在撒谎。
2 result: 他穿过衣服的牙齿穿衣服在牙齿。
3
4 source: He put on an air of innocence . 他摆出一副无辜的样子。
5 result: 他在空空空空空空空气上发现了空气。
```

Lack in capture features of different languages.

For example, Chinese tends to use less words to express the same meaning than English, but the model tends to generate the same length of sentences, so there is many repeated words in the generated sentences.

```
1 source: He ' s a friend of my brother ' s . 他是我哥哥的朋友。
2 result: 他是我弟弟弟弟弟弟的一个朋友。
```

Redundant words : Hard time finding end signal

- 1 source: You may injure yourself if you don ' t follow safety procedures . 如果你不按照安全手续来的话，你可能会受伤的。
- 2 result: 如果你分之前，你不能为止止止止止止止止止止止止止止止止止止止止止止止止止止止止止止止不能不能
- 3
- 4 source: This plane flies between Osaka and Hakodate . 这架飞机往返于大阪和函馆之间。
- 5 result: 这架飞机和大之间之间之间之间之间之间之间飞机之间之间之间之间之前坐飞机之间之间之间之间之间之

References

- <https://www.mindspore.cn/>