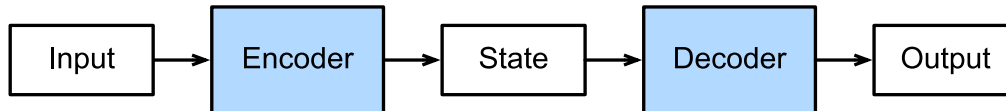


From Encoder-Decoder To Transformer

images are from <https://zh-v2.d2l.ai/>

Encoder-Decoder

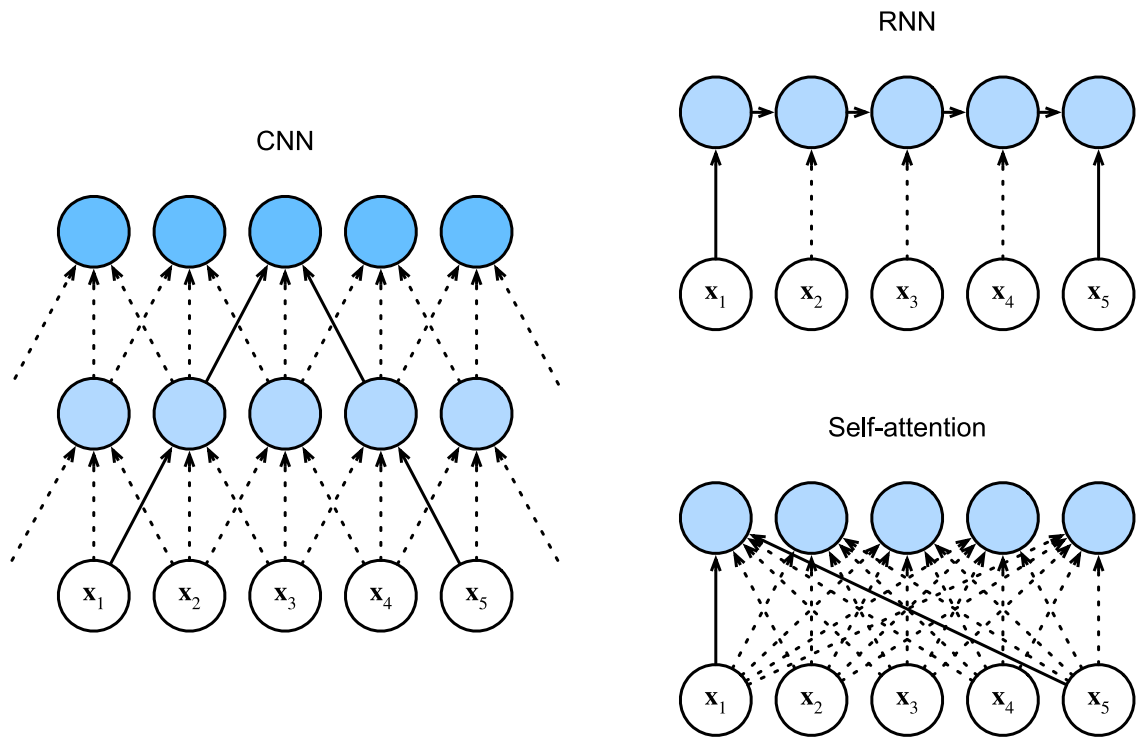


- an architecture commonly used in NLP and other types of tasks
- Encoder: take raw input and represent the input as tensors after processing (could be word2vec, neural layers, attention...)
- Decoder: mainly for outputting the result to desired form ([0, 1], probability distribution, classification, etc)

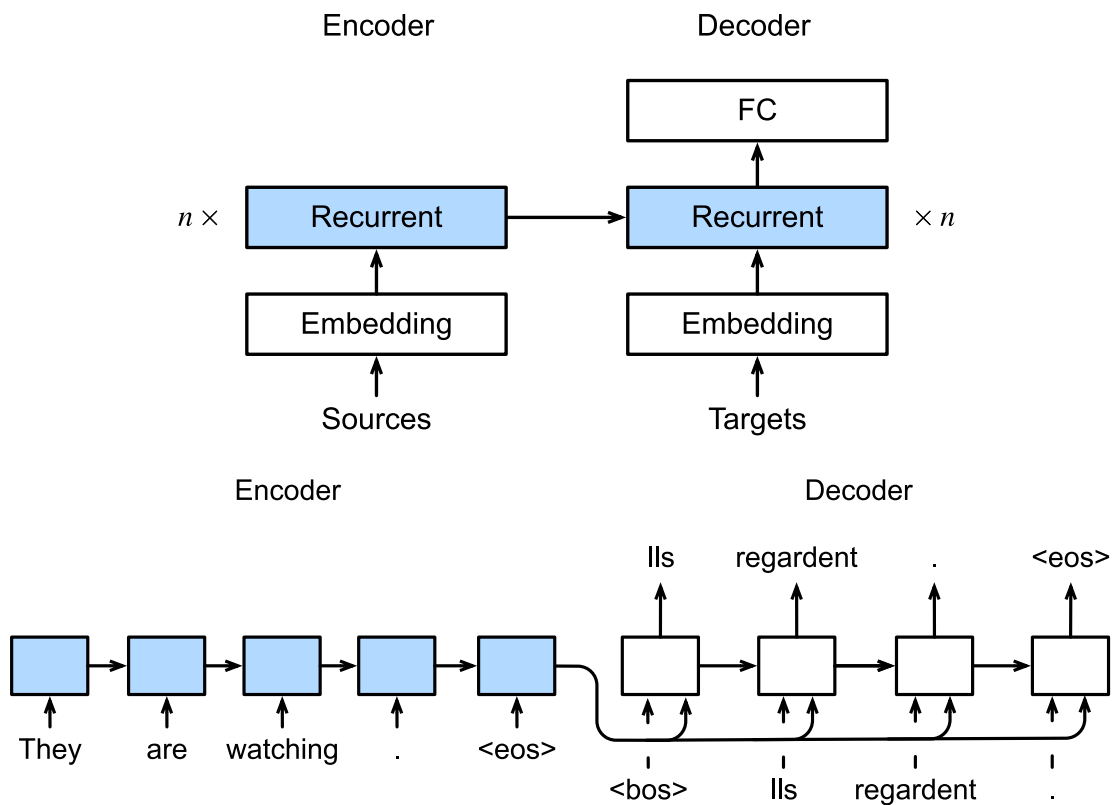
```
1  from torch import nn
2
3  class Encoder(nn.Module):
4      def __init__(self, **kwargs):
5          super(Encoder, self).__init__(**kwargs)
6
7      def forward(self, x, *args):
8          raise NotImplementedError
9
10 class Decoder(nn.Module):
11     def __init__(self, **kwargs):
12         super(Decoder, self).__init__(**kwargs)
13
14     def init_state(self, encoder_outputs, *args):
15         raise NotImplementedError
16
17     def forward(self, x, state):
18         raise NotImplementedError
```

Seq2Seq Learning

- A specific type of tasks whose input and output are both sequences of any length
- Ex. Machine Translation
- Common arch of seq2seq models:

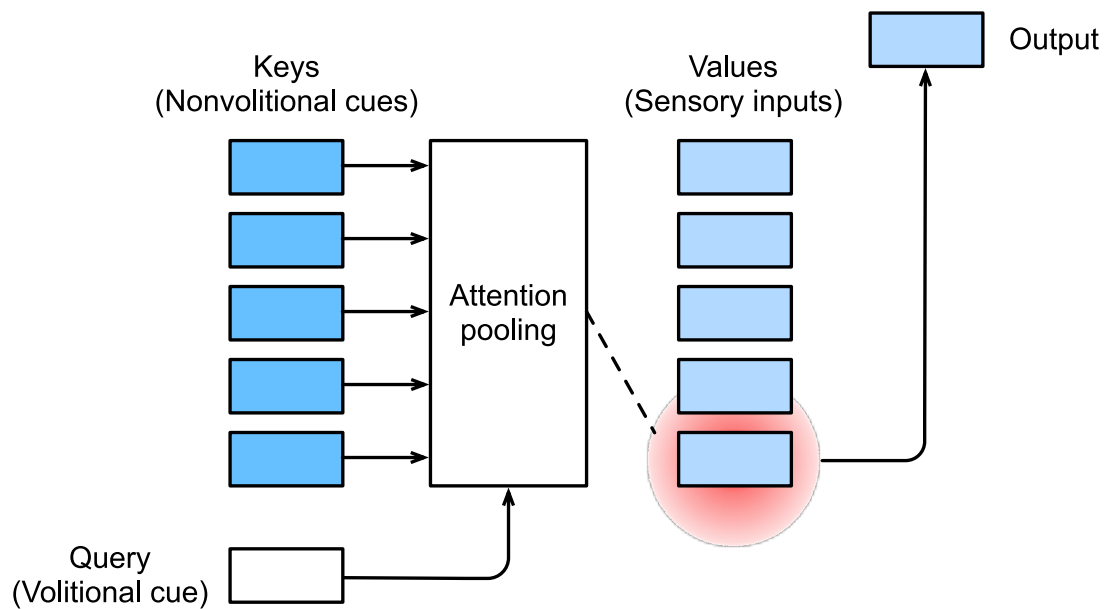


- Machine Translation using RNN



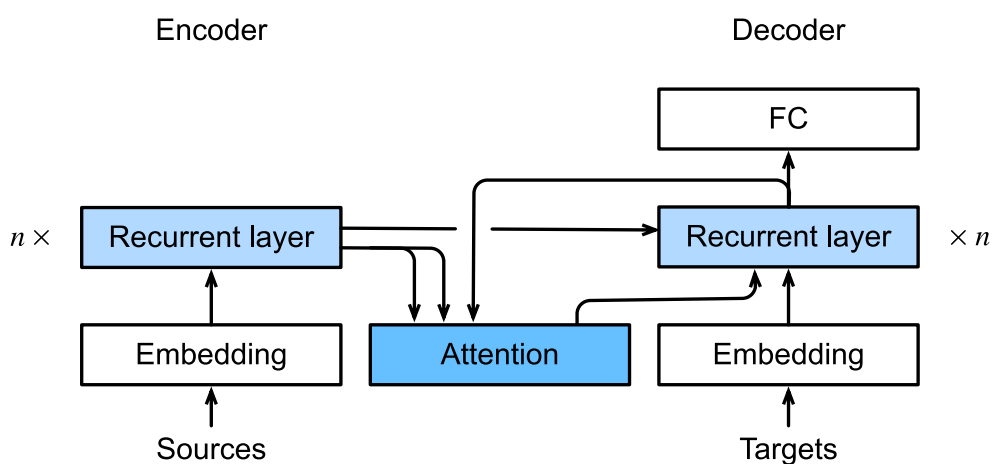
- BLEU(Bilingual Evaluation Understudy) for machine translation
 - formula: $\exp\left(\min\left(0, 1 - \frac{\text{len}_{\text{label}}}{\text{len}_{\text{pred}}}\right)\right) \prod_{n=1}^k p_n^{1/2^n}$
 - where p_n represent the **n-gram** accuracy

Attention Mechanism & Attention Score

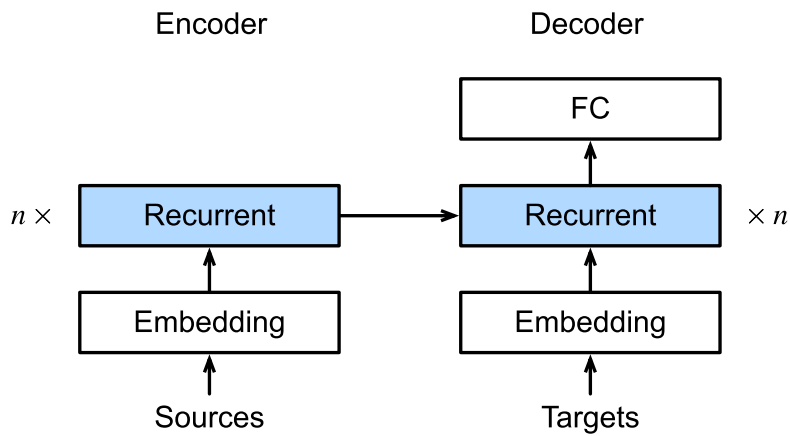


- Attention Mechanism, **KVQ**
 - **key**: what is presented
 - **value**: sensory inputs(?)
 - **Query**: what we are interested
 - The idea is to using Query to find "important" **key**s
- Attention Score, $\alpha(x, x_i)$
 - model the relationship(importance, similarity) of **keys** & **Query**s
 - Kernel Regression
 - $\alpha(x, x_i) = \frac{K(x-x_i)}{\sum_{j=1}^n K(x-x_j)}$
 - Additive Attention
 - $a(\mathbf{q}, \mathbf{k}) = \mathbf{w}_v^\top \tanh(\mathbf{W}_q \mathbf{q} + \mathbf{W}_k \mathbf{k}) \in \mathbb{R}$,
 - Scaled Dot-Product Attention
 - $a(\mathbf{q}, \mathbf{k}) = \mathbf{q}^\top \mathbf{k} / \sqrt{d}$
 - Matrix form: $\text{softmax}\left(\frac{\mathbf{QK}^\top}{\sqrt{d}}\right) \mathbf{V} \in \mathbb{R}^{n \times v}$.

Seq2Seq with Attention



- Notice the difference with



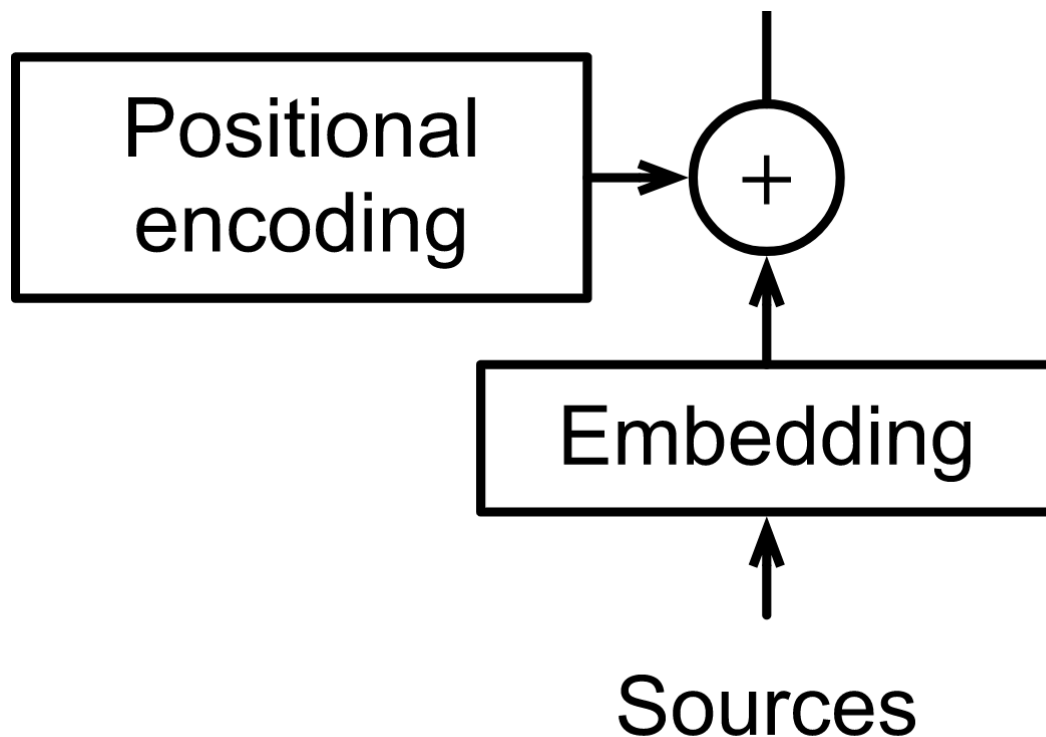
- Here, the **query** is decoder's input, **key** & **value** are both encoders output (final hidden state)

Self-attention

- Self-attention means $Queries = Values = Keys = X(input)$
- So we are trying to find the relationship between one token x_i with other tokens
- $y_i = f(x_i, (x_1, x_1), \dots, (x_n, x_n)) \in \mathbb{R}^d$, where x_i is **query** and (x_j, x_j) is **Key-Value**

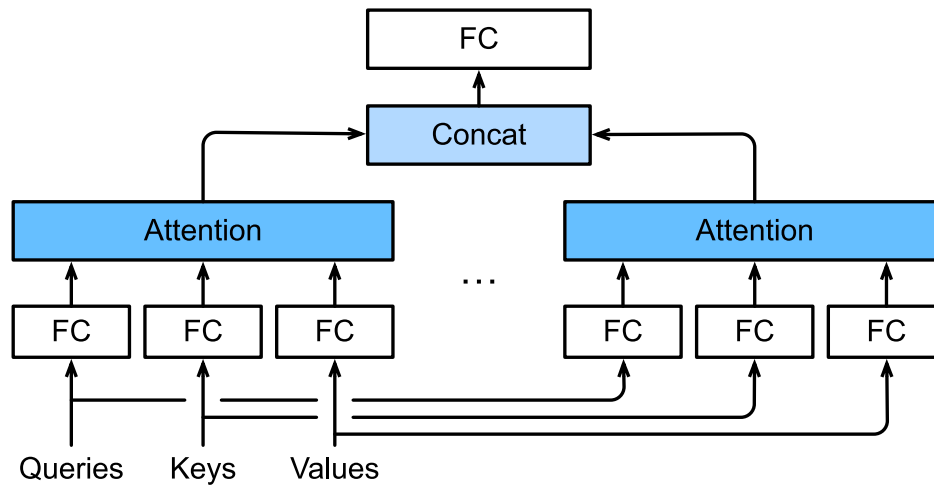
Position Encoding

- Self-attention does not contain information about relative positions (of tokens)
- Position Encoding aims to "encode" some relative position information to the input X



- A commonly used position encoding method is using these *sin* and *cos*
 - for the Position Encoding Matrix P
 - $P_{i,2j} = \sin\left(\frac{i}{10000^{2j/d}}\right)$
 - $P_{i,2j+1} = \cos\left(\frac{i}{10000^{2j/d}}\right)$

Multi-head Attention



- Multi-head Attention aims to capture different "relationships" between `Query` and `Key` using multiple parallel attention layers and concat them to get the final result.
- Mathematically:
 - $\mathbf{h}_i = f(\mathbf{W}_i^{(q)} \mathbf{q}, \mathbf{W}_i^{(k)} \mathbf{k}, \mathbf{W}_i^{(v)} \mathbf{v}) \in \mathbb{R}^{p_v}$, where f is some kind of attention function and \mathbf{h}_i is the i_{th} head
 - $result = \mathbf{W}_o \begin{bmatrix} \mathbf{h}_1 \\ \vdots \\ \mathbf{h}_h \end{bmatrix} \in \mathbb{R}^{p_o}$

```

1  class MultiHeadAttention(nn.Module):
2      def __init__(self, key_size, query_size, value_size, num_hiddens,
3                  num_heads, dropout, bias=False, **kwargs):
4          super(MultiHeadAttention, self).__init__(**kwargs)
5          self.num_heads = num_heads
6          self.attention = d2l.DotProductAttention(dropout)
7          self.w_q = nn.Linear(query_size, num_hiddens, bias=bias)
8          self.w_k = nn.Linear(key_size, num_hiddens, bias=bias)
9          self.w_v = nn.Linear(value_size, num_hiddens, bias=bias)
10         self.w_o = nn.Linear(num_hiddens, num_hiddens, bias=bias)
11
12     def forward(self, queries, keys, values, valid_lens):
13         # assuming num_queries = num_keys = num_values
14
15         # initial queries:
16         # (batch_size, num_queries, num_hiddens)
17         # transformed queries:
18         # (batch_size * num_heads, num_queries, num_hiddens/num_heads)
19         queries = transpose_qkv(self.w_q(queries), self.num_heads)
20         keys = transpose_qkv(self.w_k(keys), self.num_heads)
21         values = transpose_qkv(self.w_v(values), self.num_heads)
22
23         if valid_lens is not None:
24             valid_lens = torch.repeat_interleave(
25                 valid_lens, repeats=self.num_heads, dim=0)
26
27         # (batch_size * num_heads, num_queries, num_hiddens/num_heads)
28         output = self.attention(queries, keys, values, valid_lens)
29
30         # (batch_size, num_queries, num_hiddens)

```

```

31     output_concat = transpose_output(output, self.num_heads)
32     return self.W_o(output_concat)

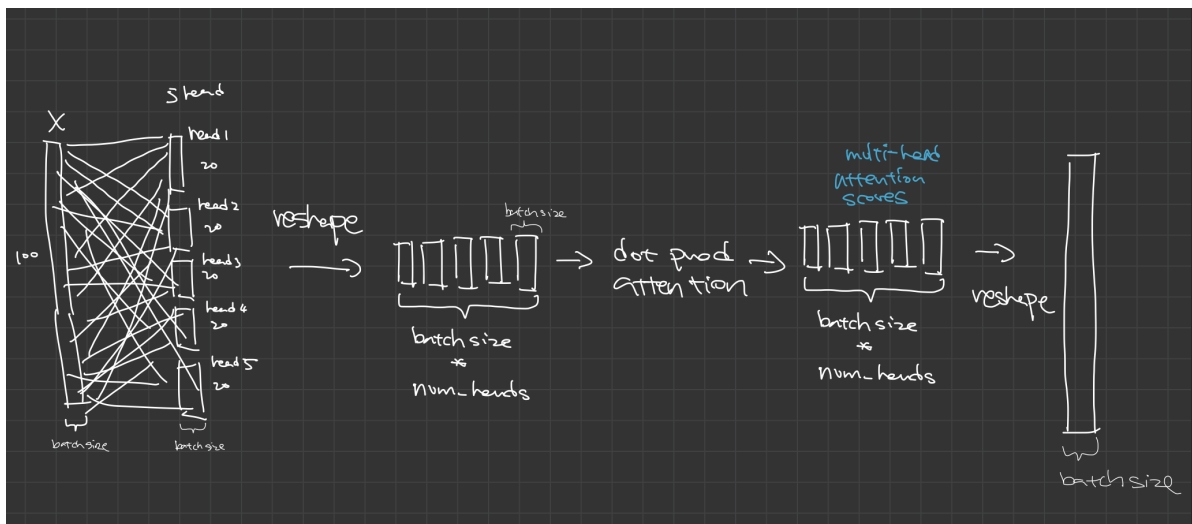
```

```

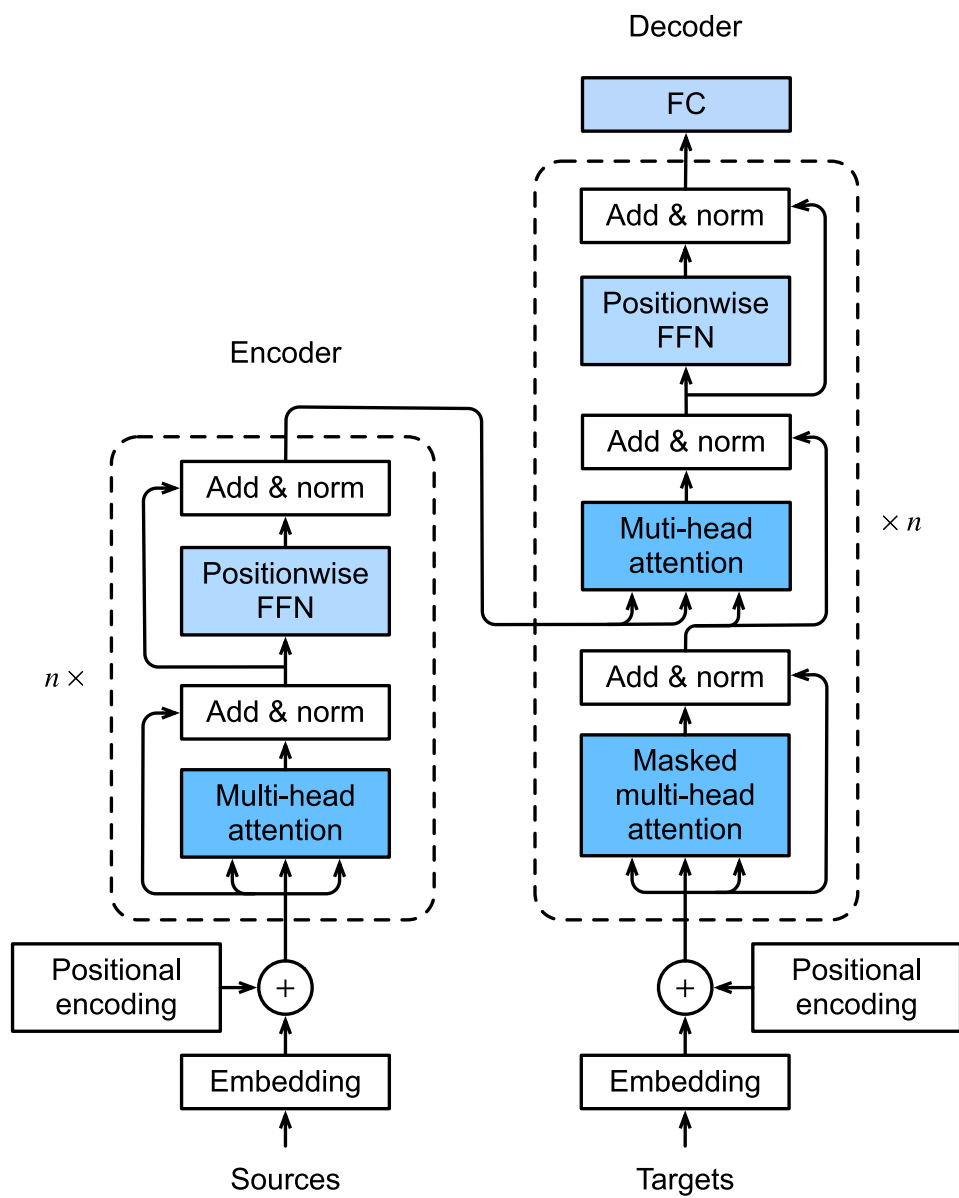
1  def transpose_qkv(X, num_heads):
2      # (batch_size, num_queries, num_hiddens)
3
4      # (batch_size, num_queries, num_heads, num_hiddens/num_heads)
5      X = X.reshape(X.shape[0], X.shape[1], num_heads, -1)
6
7      # (batch_size, num_heads, num_queries, num_hiddens/num_heads)
8      X = X.permute(0, 2, 1, 3)
9
10     # (batch_size * num_heads, num_queries, num_hiddens/num_heads)
11     return X.reshape(-1, X.shape[2], X.shape[3])
12
13
14  def transpose_output(X, num_heads):
15      """ reverse `transpose_qkv` """
16      X = X.reshape(-1, num_heads, X.shape[1], X.shape[2])
17      X = X.permute(0, 2, 1, 3)
18      return X.reshape(X.shape[0], X.shape[1], -1)

```

- Shaping



Transformer



Annotated graph

