

```

import torch
import torch.nn as nn
import math

class PositionalEncoding(nn.Module):
    def __init__(self, d_model, max_len=5000):
        super(PositionalEncoding, self).__init__()
        pe = torch.zeros(max_len, d_model)
        position = torch.arange(0, max_len, dtype=torch.float).unsqueeze(1)
        div_term = torch.exp(torch.arange(0, d_model, 2).float() * (-math.log(10000.0) / d_model))
        pe[:, 0::2] = torch.sin(position * div_term)
        pe[:, 1::2] = torch.cos(position * div_term)
        pe = pe.unsqueeze(0)
        self.register_buffer('pe', pe)

    def forward(self, x):
        return x + self.pe[:, :x.size(1)]

class MultiHeadAttention(nn.Module):
    def __init__(self, d_model, num_heads):
        super(MultiHeadAttention, self).__init__()
        assert d_model % num_heads == 0

        self.d_k = d_model // num_heads
        self.num_heads = num_heads

        self.q_linear = nn.Linear(d_model, d_model)
        self.k_linear = nn.Linear(d_model, d_model)
        self.v_linear = nn.Linear(d_model, d_model)
        self.out = nn.Linear(d_model, d_model)

    def forward(self, q, k, v, mask=None):
        batch_size = q.size(0)

        def transform(x, linear):
            x = linear(x)
            x = x.view(batch_size, -1, self.num_heads, self.d_k)
            return x.transpose(1, 2)

        q = transform(q, self.q_linear)
        k = transform(k, self.k_linear)
        v = transform(v, self.v_linear)

        scores = torch.matmul(q, k.transpose(-2, -1)) / math.sqrt(self.d_k)

```

```

if mask is not None:
    scores = scores.masked_fill(mask == 0, -1e9)

    attn = torch.softmax(scores, dim=-1)
    output = torch.matmul(attn, v)

    output = output.transpose(1, 2).contiguous()
    output = output.view(batch_size, -1, self.num_heads * self.d_k)

    return self.out(output)

```

```

class FeedForward(nn.Module):

```

```

    def __init__(self, d_model, d_ff):
        super(FeedForward, self).__init__()
        self.linear1 = nn.Linear(d_model, d_ff)
        self.linear2 = nn.Linear(d_ff, d_model)

    def forward(self, x):
        return self.linear2(torch.relu(self.linear1(x)))

```

```

class TransformerBlock(nn.Module):

```

```

    def __init__(self, d_model, num_heads, d_ff, dropout=0.1):
        super(TransformerBlock, self).__init__()
        self.attention = MultiHeadAttention(d_model, num_heads)
        self.feed_forward = FeedForward(d_model, d_ff)
        self.norm1 = nn.LayerNorm(d_model)
        self.norm2 = nn.LayerNorm(d_model)
        self.dropout = nn.Dropout(dropout)

    def forward(self, x, mask=None):
        attn_out = self.attention(x, x, x, mask)
        x = self.norm1(x + self.dropout(attn_out))
        ff_out = self.feed_forward(x)
        x = self.norm2(x + self.dropout(ff_out))
        return x

```

```

class TransformerEncoder(nn.Module):

```

```

    def __init__(self, vocab_size, d_model, num_layers, num_heads, d_ff, max_len=512):
        super(TransformerEncoder, self).__init__()
        self.embedding = nn.Embedding(vocab_size, d_model)
        self.positional_encoding = PositionalEncoding(d_model, max_len)
        self.layers = nn.ModuleList([
            TransformerBlock(d_model, num_heads, d_ff) for _ in range(num_layers)
        ])
        self.dropout = nn.Dropout(0.1)

```

```
def forward(self, src, mask=None):
    x = self.embedding(src)
    x = self.positional_encoding(x)
    x = self.dropout(x)

    for layer in self.layers:
        x = layer(x, mask)

    return x

if __name__ == "__main__":
    vocab_size = 10000
    d_model = 512
    num_layers = 6
    num_heads = 8
    d_ff = 2048

    model = TransformerEncoder(vocab_size, d_model, num_layers, num_heads, d_ff)
    dummy_input = torch.randint(0, vocab_size, (32, 100))
    output = model(dummy_input)
    print(output.shape)
```