```
from google.colab import drive
drive.mount('/content/drive')

→ Mounted at /content/drive
```

```
+ Code - + Text
```

# Dependencies

```
!pip install -U spacy==3.6.0
!python -m spacy download en_core_web_sm
!python -m spacy download de_core_news_sm
!pip install torchdata
!pip install -U torchtext
!pip install portalocker>=2.0.0
!pip install seaborn
      Downloading nvidia_cufft_cu12-11.0.2.54-py3-none-manylinux1_x86_64.whl (121.6 MB)
                                                  - 121.6/121.6 MB 14.5 MB/s eta 0:00:00
    Collecting nvidia-curand-cu12==10.3.2.106 (from torch==2.1.1->torchtext)
      Downloading nvidia_curand_cu12-10.3.2.106-py3-none-manylinux1_x86_64.whl (56.5 MB)
                                                  - 56.5/56.5 MB 30.0 MB/s eta 0:00:00
    Collecting nvidia-cusolver-cu12==11.4.5.107 (from torch==2.1.1->torchtext)
      Downloading nvidia_cusolver_cu12-11.4.5.107-py3-none-manylinux1_x86_64.whl (124.2 MB)
                                                   124.2/124.2 MB 7.7 MB/s eta 0:00:00
    Collecting nvidia-cusparse-cu12==12.1.0.106 (from torch==2.1.1->torchtext)
      Downloading nvidia_cusparse_cu12-12.1.0.106-py3-none-manylinux1_x86_64.whl (196.0 MB)
                                                   196.0/196.0 MB 4.0 MB/s eta 0:00:00
    Collecting nvidia-nccl-cu12==2.18.1 (from torch==2.1.1->torchtext)
      Downloading nvidia_nccl_cu12-2.18.1-py3-none-manylinux1_x86_64.whl (209.8 MB)
                                                  - 209.8/209.8 MB 4.5 MB/s eta 0:00:00
    Collecting nvidia-nvtx-cu12==12.1.105 (from torch==2.1.1->torchtext)
      Downloading nvidia_nvtx_cu12-12.1.105-py3-none-manylinux1_x86_64.whl (99 kB)
                                                   99.1/99.1 kB 13.1 MB/s eta 0:00:00
    Requirement already satisfied: triton==2.1.0 in /usr/local/lib/python3.10/dist-packages (from torch==2.1.1->torchtext) (2.1.0)
    Requirement already satisfied: urllib3>=1.25 in /usr/local/lib/python3.10/dist-packages (from torchdata==0.7.1->torchtext) (2.0.7)
    Collecting nvidia-nvjitlink-cu12 (from nvidia-cusolver-cu12==11.4.5.107->torch==2.1.1->torchtext)
      Downloading nvidia_nvjitlink_cu12-12.3.101-py3-none-manylinux1_x86_64.whl (20.5 MB)
                                                   20.5/20.5 MB 90.4 MB/s eta 0:00:00
    Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests->torchtext) (3.3.2)
    Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests->torchtext) (3.4)
    Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests->torchtext) (2023.7.22)
    Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.10/dist-packages (from jinja2->torch==2.1.1->torchtext) (2.1.3)
    Requirement already satisfied: mpmath>=0.19 in /usr/local/lib/python3.10/dist-packages (from sympy->torch==2.1.1->torchtext) (1.3.0)
    Installing collected packages: nvidia-nvtx-cu12, nvidia-nvjitlink-cu12, nvidia-nccl-cu12, nvidia-curand-cu12, nvidia-cufft-cu12, nvidia-c
      Attempting uninstall: torch
        Found existing installation: torch 2.1.0+cu118
        Uninstalling torch-2.1.0+cu118:
          Successfully uninstalled torch-2.1.0+cu118
      Attempting uninstall: torchdata
        Found existing installation: torchdata 0.7.0
        Uninstalling torchdata-0.7.0:
          Successfully uninstalled torchdata-0.7.0
      Attempting uninstall: torchtext
        Found existing installation: torchtext 0.16.0
        Uninstalling torchtext-0.16.0:
          Successfully uninstalled torchtext-0.16.0
    ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviour is the source o
    torchaudio 2.1.0+cu118 requires torch==2.1.0, but you have torch 2.1.1 which is incompatible.
    torchvision 0.16.0+cu118 requires torch==2.1.0, but you have torch 2.1.1 which is incompatible.
    Successfully installed nvidia-cublas-cu12-12.1.3.1 nvidia-cuda-cupti-cu12-12.1.105 nvidia-cuda-nvrtc-cu12-12.1.105 nvidia-cuda-runtime-cu
    Requirement already satisfied: seaborn in /usr/local/lib/python3.10/dist-packages (0.12.2)
    Requirement already satisfied: numpy!=1.24.0,>=1.17 in /usr/local/lib/python3.10/dist-packages (from seaborn) (1.23.5) Requirement already satisfied: pandas>=0.25 in /usr/local/lib/python3.10/dist-packages (from seaborn) (1.5.3)
    Requirement already satisfied: matplotlib!=3.6.1,>=3.1 in /usr/local/lib/python3.10/dist-packages (from seaborn) (3.7.1)
    Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (1.2.0
    Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (0.12.1)
    Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (4.44
    Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (1.4.
    Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (23.2)
    Requirement already satisfied: pillow>=6.2.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (9.4.0)
    Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (3.1.1
    Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (2
    Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas>=0.25->seaborn) (2023.3.post1)
    Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.7->matplotlib!=3.6.1,>=3.1->s
```

# Transformer Assignment

In this assignment, you will be trying your hand at understanding transformers, their architecture, and their difference in-terms of basic RNNs. The assignment is divided in 2 sections.

Section 1:

You will be implmenting a basic RNN cell, RNN Class and an RNN Classifier

Section 2:

You will be implementing a Transformer based Text classifier using components such as Multi-head Attention Module, Positional Encoding Module and Encoder

· Section 3:

In order to experiment with Decoders for a Transformer, we will be implementing a Transformer Based Machine Translation class using modules of Section 2, a Decoder, Attention Masks and Seq-Seq Module

```
import math
import torch
import time
import numpy as np
import torch.nn as nn
import torch.nn.functional as F
import pandas as pd
import seaborn as sns
from torch.utils.data import DataLoader
from torchtext.datasets import AG_NEWS
from torch.utils.data.dataset import random_split
from torchtext.data.functional import to_map_style_dataset
from torchtext.data.utils import get_tokenizer
from torchtext.vocab import build_vocab_from_iterator
from torchtext.datasets import Multi30k
from typing import Iterable, List
```

# Section 1: Recurrent Neural Networks (RNN)

Each RNN Cell should contain 2 components: an Input Unit and a Hidden Unit. The Hidden state is the part of the RNN that remembers context about previous data present in the sequence. The current time step's hidden state is calculated using information of the previous time steps hidden state and the current input. This process helps to retain information on what the model saw in the previous time step when processing the current time steps information.

RNNs will look and function as follows,



The hidden state any given time t is given by,

$$input_t = (x_t \cdot W_x^t + b_x^t)$$
  
 $prev\_state = (h_{t-1} \cdot W_h^t + b_h^t)$   
 $h_t = tanh(input_t + prev\_state)$ 

The output at any give time t is given by,

$$y_t = h_t \cdot W_y^t + b_y$$

Note: All the connections in RNN have weights and biases.

Your job is to implement the formulae above.

# 1.1 A Single RNN Cell

```
class RNNCell(torch.nn.Module):
   RNNCell is a single cell that takes x_t and h_{t_1} as input and outputs h_t.
   def __ı
       _init__(self, input_dim: int, hidden_dim: int):
      Constructor of RNNCell.
      Inputs:
      - input dim: Dimension of the input x t
      - hidden_dim: Dimension of the hidden state h_{t-1} and h_t
      .....
      # We always need to do this step to properly implement the constructor
      super(RNNCell, self).__init__()
      self.linear_x, self.linear_h, self.non_linear = None, None, None
      # 1. Define the linear transformation layers for the attributes
          (set to None above) to correspond to the W_x and W_h in the formulae.
          Remember to include bias in the linear layers.
          (Refer to nn.Linear documentation https://pytorch.org/docs/stable/generated/torch.nn.Linear.html)
      # 2. Define the non_linear layer. (You can use tanh as describe bove).
      # Linear transformations
      self.linear_x = nn.Linear(input_dim, hidden_dim) # bias is included by default
      self.linear_h = nn.Linear(hidden_dim, hidden_dim)
      # Non-linear activation function
      self.non_linear = nn.Tanh()
      END OF YOUR CODE
      def forward(self, x_cur: torch.Tensor, h_prev: torch.Tensor):
      Compute h_t given x_t and h_{t-1}.
      Inputs:
      - x_cur: x_t, a tensor with the same of BxC, where B is the batch size and
       C is the channel dimension.
      - h_prev: h_{t-1}, a tensor with the same of BxH, where H is the channel
       dimension.
      h_{cur} = None
      # TODO: Run the linear transformation layers to compute x_t and consume
      # go non-linear layer.
      combined = self.linear_x(x_cur) + self.linear_h(h_prev)
      h_cur = self.non_linear(combined)
      END OF YOUR CODE
      return h_cur
# Let's run a sanity check of your model
x = torch.randn((2, 8)) # Input Dim
h = torch.randn((2, 16)) # Hidden Dim
model = RNNCell(8, 16)
y = model(x, h)
assert len(y.shape) == 2 and y.shape[0] == 2 and y.shape[1] == 16
print(y.shape)
   torch.Size([2, 16])
```

```
RNN is a single-layer (stack) RNN by connecting multiple RNNCell together in a single
   direction, where the input sequence is processed from left to right.
   def __init__(self, input_dim: int, hidden_dim: int):
     Constructor of the RNN module.
     Inputs:
     - input_dim: Dimension of the input x_t
     - hidden_dim: Dimension of the hidden state h_{t-1} and h_t
     super(RNN, self).__init__()
     self.hidden_dim = hidden_dim
     # TODO: Define the RNNCell.
     self.rnn_cell = RNNCell(input_dim, hidden_dim)
     END OF YOUR CODE
     def forward(self, x: torch.Tensor):
     Compute the hidden representations for every token in the input sequence.
     Input:
     - x: A tensor with the shape of BxLxC, where B is the batch size, L is the squence
       length, and C is the channel dimmension
     Return:
      - h: A tensor with the shape of BxLxH, where H is the hidden dimension of RNNCell
     b = x.shape[0]
     seq_len = x.shape[1]
     # initialize the hidden dimension
     init_h = x.new_zeros((b, self.hidden_dim))
     h = None
     # TODO: Compute the hidden representation for every token in the input
     # from left to right as per the formula stated above
     h_stack = []
     h_prev = init_h
     for t in range(seq_len):
        h_t = self.rnn_cell(x[:, t, :], h_prev)
        # Store the current hidden state
        h_stack.append(h_t)
        h_prev = h_t
     # Convert the list of hidden states to a tensor
     h = torch.stack(h_stack, dim=1)
     END OF YOUR CODE
     return h
# Let's run a sanity check of your model
x = torch.randn((2, 10, 8))
model = RNN(8, 16)
y = model(x)
assert len(y.shape) == 3
for dim, dim_gt in zip(y.shape, [2, 10, 16]):
   assert dim == dim_gt
print(y.shape)
   torch.Size([2, 10, 16])
```

class RNN(torch.nn.Module):

```
class RNNClassifier(nn.Module):
  A RNN-based classifier for text classification. It first converts tokens into word embeddings.
  And then feeds the embeddings into a RNN, where the hidden representations of all tokens are
  then averaged to get a single embedding of the sentence. It will be used as input to a linear
  classifier.
  def __init__(self,
        vocab_size: int, embed_dim: int, rnn_hidden_dim: int, num_class: int, pad_token: int
     ):
     Constructor.
     Inputs:
     - vocab_size: Vocabulary size, indicating how many tokens we have in total.
     - embed_dim: The dimension of word embeddings
     - rnn_hidden_dim: The hidden dimension of the RNN.
     - num_class: Number of classes.
     - pad_token: The index of the padding token.
     super(RNNClassifier, self).__init__()
     # word embedding layer
     self.embedding = nn.Embedding(vocab_size, embed_dim, padding_idx=pad_token)
     self.rnn, self.fc = None, None
     # TODO: Define the RNN and the classification layer.
     # Define the RNN
     self.rnn = RNN(embed_dim, rnn_hidden_dim)
     # Define the classification layer
     self.fc = nn.Linear(rnn_hidden_dim, num_class)
     END OF YOUR CODE
     def init_weights(self):
     initrange = 0.5
     self.embedding.weight.data.uniform_(-initrange, initrange)
     self.fc.weight.data.uniform_(-initrange, initrange)
     self.fc.bias.data.zero_()
  def forward(self, text):
     Get classification scores (logits) of the input.
     Input:
      - text: Tensor with the shape of BxLxC.
     Return:
     - logits: Tensor with the shape of BxK, where K is the number of classes
     # get word embeddings
     embedded = self.embedding(text)
     logits = None
     # TODO: Compute logits of the input.
     # Get word embeddings
     embedded = self.embedding(text)
     # Pass embeddings through RNN
     rnn_output = self.rnn(embedded)
     # Average the outputs of RNN across the sequence length dimension
     avg_output = rnn_output.mean(dim=1)
     # Compute logits
     logits = self.fc(avg_output)
     END OF YOUR CODE
```

h\_tracker = {}

```
# Sanity check!!!
vocab_size = 10
embed_dim = 16
rnn_hidden_dim = 32
num_class = 4

x = torch.arange(vocab_size).view(1, -1)
x = torch.cat((x, x), dim=0)
print('x.shape: {}'.format(x.shape))
model = RNNClassifier(vocab_size, embed_dim , rnn_hidden_dim, num_class, 0)
y = model(x)
assert len(y.shape) == 2 and y.shape[0] == 2 and y.shape[1] == num_class
print(y.shape)

x.shape: torch.Size([2, 10])
torch.Size([2, 4])
```

# → Data Loader

```
# check here for details https://github.com/pytorch/text/blob/main/torchtext/data/utils.py#L52-#L166
from torchtext.data.utils import get_tokenizer
# check here for details https://github.com/pytorch/text/blob/main/torchtext/vocab/vocab_factory.py#L65-L113
from torchtext.vocab import build_vocab_from_iterator
# Documentation of DataLoader https://pytorch.org/docs/stable/data.html#torch.utils.data.DataLoader
from torch.utils.data import DataLoader
# A tokenizer splits a input setence into a set of tokens, including those puncuation
# For example
# >>> tokens = tokenizer("You can now install TorchText using pip!")
# >>> tokens
# >>> ['you', 'can', 'now', 'install', 'torchtext', 'using', 'pip', '!']
tokenizer = get_tokenizer('basic_english')
train_iter = AG_NEWS(split='train')
def yield_tokens(data_iter):
   for _, text in data_iter:
       yield tokenizer(text)
# Creates a vocab object which maps tokens to indices
# Check here for details https://github.com/pytorch/text/blob/main/torchtext/vocab/vocab.py
vocab = build_vocab_from_iterator(yield_tokens(train_iter), specials=["<unk>"])
# The specified token will be returned when a out-of-vocabulary token is queried.
vocab.set_default_index(vocab["<unk>"])
text_pipeline = lambda x: vocab(tokenizer(x))
label_pipeline = lambda x: int(x) - 1
# The padding token we need to use
# The returned indices are always in an array
PAD_TOKEN = vocab(tokenizer('<pad>'))
assert len(PAD_TOKEN) == 1
PAD_TOKEN = PAD_TOKEN[0]
# Merges a list of samples to form a mini-batch of Tensor(s)
def collate_batch(batch):
   Input:
   - batch: A list of data in a mini batch, where the length denotes the batch size.
     The actual context depends on a particular dataset. In our case, each position
     contains a label and a Tensor (tokens in a sentence).
   Returns:
   - batched_label: A Tensor with the shape of (B,)
   - batched_text: A Tensor with the shape of (B, L, C), where L is the sequence length
     and C is the channeld dimension
   label_list, text_list, text_len_list = [], [], []
   for (_label, _text) in batch:
       label_list.append(label_pipeline(_label))
       processed_text = torch.tensor(text_pipeline(_text), dtype=torch.int64)
       text_list.append(processed_text)
       text_len_list.append(processed_text.size(0))
   batched_label, batched_text = None, None
   \# TODO: Pad the text tensor in the mini batch so that they have the same \#
   # length. Specifically, you need to calculate the maximum length in the
   # batch and then add the token PAD_TOKEN to the end of those
   # shorter sentences. (Try printing a few data points to understand why)
   # Find the maximum length in this batch
   max_len = max(text_len_list)
   # Pad each text tensor to the maximum length
   padded_texts = [torch.cat([text, torch.full((max_len - len(text),), PAD_TOKEN, dtype=torch.int64)])
                  for text in text_list]
   # Stack all the padded texts and labels to create batch tensors
   batched_label = torch.tensor(label_list, dtype=torch.int64)
   batched_text = torch.stack(padded_texts)
   END OF YOUR CODE
   return batched_label.long(), batched_text.long()
# Now, let's check what the batched data looks like
train_iter = AG_NEWS(split='train')
```

```
dataloader = DataLoader(train_iter, batch_size=8, shuffle=False, collate_fn=collate_batch)
for idx, (label, data) in enumerate(dataloader):
    if idx > 0:
        break
    print('label.shape: {}'.format(label.shape))
    print('label: {}'.format(label))
    print('data.shape: {}'.format(data.shape))

    label.shape: torch.Size([8])
    label: tensor([2, 2, 2, 2, 2, 2, 2])
    data.shape: torch.Size([8, 49])

labels=set()
labels.update([entry[0] for entry in AG_NEWS(root="data")[0]])
print(labels)
    {1, 2, 3, 4}
```

# ▼ 1.4 Train & Evaluate Module

```
# logits_tracker = {}
def train(model, dataloader, loss_func, device, grad_norm_clip, optimizer):
  model.train()
  total_acc, total_count = 0, 0
  log_interval = 500
  start_time = time.time()
  global logits_tracker
  for idx, (label, text) in enumerate(dataloader):
     label = label.to(device)
     text = text.to(device)
     optimizer.zero_grad()
     logits = None
     # TODO: compute the logits of the input, get the loss, and do the
     # gradient backpropagation.
     logits = model(text)
     loss = loss_func(logits, label)
     loss.backward()
     END OF YOUR CODE
     torch.nn.utils.clip_grad_norm_(model.parameters(), grad_norm_clip)
     optimizer.step()
     total_acc += (logits.argmax(1) == label).sum().item()
     total_count += label.size(0)
     if idx % \log_{interval} == 0 and idx > 0:
        elapsed = time.time() - start_time
        print('| epoch {:3d} | {:5d}/{:5d} batches '
             '| accuracy {:8.3f}'.format(epoch, idx, len(dataloader),
                                total_acc/total_count))
        total_acc, total_count = 0, 0
        start_time = time.time()
def evaluate(model, dataloader, loss_func, device):
  model.eval()
  total_acc, total_count = 0, 0
  with torch.no_grad():
     for idx, (label, text) in enumerate(dataloader):
        label = label.to(device)
        text = text.to(device)
        # TODO: compute the logits of the input, get the loss.
        # Compute logits
        logits = model(text)
        END OF YOUR CODE
        total_acc += (logits.argmax(1) == label).sum().item()
        total_count += label.size(0)
  return total_acc/total_count
```

```
assert torch.cuda.is_available(), "Please connect to the GPU instance if working on Colab or configure the environment for Torch using GPU (C
# device = 'cuda'
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
# Hyper parameters
epochs = 3 # epoch
lr =0.0005 # learning rate
batch_size = 64 # batch size for training
word embed dim = 64
rnn_hidden_dim = 96
train_iter = AG_NEWS(split='train')
num_class = len(set([label for (label, text) in train_iter]))
vocab_size = len(vocab)
model, loss_func = None, None
# TODO: Define the classifier and loss function.
model = RNNClassifier(vocab_size, word_embed_dim, rnn_hidden_dim, num_class, PAD_TOKEN)
loss_func = nn.CrossEntropyLoss()
END OF YOUR CODE
# copy the model to the specified device (GPU)
model = model.to(device)
optimizer = torch.optim.AdamW(model.parameters(), lr=lr)
scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(optimizer, epochs, 1e-8)
total_accu = None
train_iter, test_iter = AG_NEWS()
train_dataset = to_map_style_dataset(train_iter)
test_dataset = to_map_style_dataset(test_iter)
num_train = int(len(train_dataset) * 0.95)
split_train_, split_valid_ = random_split(
   train_dataset,
   [num_train, len(train_dataset) - num_train]
train_dataloader = DataLoader(
   split_train_, batch_size=batch_size,
   shuffle=True, collate_fn=collate_batch
)
valid dataloader = DataLoader(
   split_valid_, batch_size=batch_size,
   shuffle=False, collate_fn=collate_batch
test_dataloader = DataLoader(
   test_dataset, batch_size=batch_size,
   shuffle=False, collate_fn=collate_batch
split_train_[21]
     'Gaza Strip bloodshed continues Violent clashes are continuing in the Gaza Strip on the fifth day of Israeli attacks in the region. At
    least three Palestinians were killed in early morning raids launched shortly after Israel #39;s army chief ')
```

```
# You should be able get a validation accuracy around 86%
for epoch in range(1, epochs + 1):
    # global logits_tracker
    # logits_tracker[epoch] = None
    epoch_start_time = time.time()
    train(model, train_dataloader, loss_func, device, 1, optimizer)
    accu_val = evaluate(model, valid_dataloader, loss_func, device)
    if total_accu is not None and total_accu > accu_val:
        scheduler.step()
    else:
        total_accu = accu_val
    print('-' * 59)
    print('| end of epoch {:3d} | time: {:5.2f}s |
          'valid accuracy {:8.3f} '.format(epoch,
                                            time.time() - epoch_start_time,
                                            accu_val))
    print('-' * 59)
      epoch
              1 |
                     500/ 1782 batches |
                                                      0.462
                                         accuracy
                    1000/ 1782 batches |
                                                      0.619
      epoch
               1
                                         accuracy
                    1500/ 1782 batches | accuracy
      end of epoch
                     1 | time: 85.43s | valid accuracy
                                                            0.747
      epoch
                     500/ 1782 batches |
                    1000/ 1782 batches
                                                      0.823
                                         accuracy
      epoch
      epoch
                    1500/ 1782 batches | accuracy
      end of epoch
                      2 | time: 84.34s | valid accuracy
                                                            0.851
                     500/ 1782 batches |
                                                      0.871
      epoch
                                         accuracy
      epoch
                    1000/ 1782 batches
                                                      0.878
                    1500/ 1782 batches |
      epoch
                                         accuracy
      end of epoch
                      3 | time: 84.96s | valid accuracy
                                                            0.867
```

# Section 2: Transformers - 'Attention is All you Need' : Classifier

Transformers are a type of deep learning architecture that has had a profound impact on a wide range of natural language processing (NLP) tasks and other sequence-to-sequence tasks. They are known for their ability to model long-range dependencies and their parallelization capabilities. A typical transformer model consists of several key components:



Transformers have revolutionized NLP and have been adapted for a wide range of applications beyond text, including image generation, recommendation systems, and more. For this section, we will be implementing, all but two important components, Attention-Masks and Decoder Module. They will be implemented in-depth in Section 3.

### 2.1 Multihead Attention



Multi-Head Attention can be mathematically explained as follows:

Let's assume we have a sequence of input vectors  $(X = x_1, x_2, \dots, x_n)$ , where  $(x_i)$  represents the (i)-th element of the sequence. Each  $x_i$  is typically a vector, such as a word embedding in natural language processing.

#### 1. Single Attention Head:

• In a single attention head, we compute attention scores  $(A_{ij})$  between every pair of input elements  $(x_i)$  and  $(x_j)$ . These scores are computed using a compatibility function, often a dot product or a learned linear transformation followed by a softmax activation:

 $A_{ij} = softmax(\frac{(Q_i x_i)^T (K_j x_j)}{\sqrt{d_k}})$  Where  $Q_i$  and  $K_j$  are learned linear transformations of the input vectors  $x_i$  and  $x_j$ , and  $d_k$  is the dimension of the key vectors.

o The attention scores are used to compute weighted representations of the input sequence:

Attention(X) = 
$$\sum_{i=1}^{n} A_{ij} V_j$$

Where  $V_i$  is a learned linear transformation of the input vector  $x_i$ .

#### 2. Multiple Attention Heads:

 $\circ$  In Multi-Head Attention, we use H attention heads in parallel. Each head has its own sets of learned parameters for Q, K, and V, resulting in H sets of attention scores and weighted representations.

$$\textit{MultiHead}(X) = \textit{Concatenate}(\textit{Head}_1\,, \textit{Head}_2\,, \dots\,, \textit{Head}_H\,).\,\textit{W}^{\,\,0}$$

Where  $W^{O}$  is another learned linear transformation applied to the concatenated outputs, and  $Head_{i}$  represents the output of the i-th attention head.

```
class MultiHeadAttention(nn.Module):
   A module that computes multi-head attention given query, key, and value tensors.
       _init__(self, input_dim: int, num_heads: int):
      Constructor.
      Inputs:
      - input_dim: Dimension of the input query, key, and value. Here we assume they all have
       the same dimensions. But they could have different dimensions in other problems.
      - num_heads: Number of attention heads
      super(MultiHeadAttention, self).__init__()
      assert input_dim % num_heads == 0 # Check if we can get back the original Dimensions!
      self.input_dim = input_dim
      self.num_heads = num_heads
      self.dim_per_head = input_dim // num_heads
      # TODO: Define the linear transformation layers for key, value, and query.#
      # Also define the output layer.
      self.query_layer = nn.Linear(input_dim, input_dim)
      self.key_layer = nn.Linear(input_dim, input_dim)
      self.value_layer = nn.Linear(input_dim, input_dim)
      self.out_layer = nn.Linear(input_dim, input_dim)
      END OF YOUR CODE
      self.scores = None
   def forward(self, query: torch.Tensor, key: torch.Tensor, value: torch.Tensor, mask: torch.Tensor=None):
      Compute the attended feature representations.
      Inputs:
      - query: Tensor of the shape BxLxC, where B is the batch size, L is the sequence length,
       and C is the channel dimension
      - key: Tensor of the shape BxLxC
      - value: Tensor of the shape BxLxC
      - mask: Tensor indicating where the attention should *not* be performed
      b = query.shape[0]
      dot_prod_scores = None
      # TODO: Compute the scores based on dot product between transformed query,#
      # key, and value. You may find torch.matmul helpful, whose documentation #
      # can be found at
      # https://pytorch.org/docs/stable/generated/torch.matmul.html#torch.matmul#
      # Remember to devide the doct product similarity scores by square root of #
      # the channel dimension per head.
      # Since no for loops are allowed here, think of how to use tensor reshape #
      # to process multiple attention heads at the same time.
      # Transform query, key, value
      query_t = self.query_layer(query).view(b, -1, self.num_heads, self.dim_per_head).transpose(1, 2)
      key_t = self.key_layer(key).view(b, -1, self.num_heads, self.dim_per_head).transpose(1, 2)
      value_t = self.value_layer(value).view(b, -1, self.num_heads, self.dim_per_head).transpose(1, 2)
      # Dot product between query and key
      dot_prod_scores = torch.matmul(query_t, key_t.transpose(-2, -1)) / (self.dim_per_head ** 0.5)
      END OF YOUR CODE
      if mask is not None:
         # We simply set the similarity scores to be near zero for the positions
         # where the attention should not be done. Think of why we do this.
         dot_prod_scores = dot_prod_scores.masked_fill(mask == 0, -1e9)
      out = None
      # TODO: Compute the attention scores, which are then used to modulate the #
      # value tensor. Finally concate the attended tensors from multiple heads #
      # and feed it into the output layer. You may still find torch.matmul
```

```
# helpful.
                                                                #
      # Again, think of how to use reshaping tensor to do the concatenation.
      # Compute the attention scores
      attention_scores = F.softmax(dot_prod_scores, dim=-1)
      # Apply attention scores to the value tensor
      attended_values = torch.matmul(attention_scores, value_t)
      # Concatenate and pass through the final linear layer
      attended\_values = attended\_values.transpose(1, 2).contiguous().view(b, -1, self.input\_dim)
      out = self.out_layer(attended_values)
      END OF YOUR CODE
      return out
# Sanity Check
x = torch.randn((2, 10, 8))
mask = torch.randn((2.10)) > 0.5
mask = mask.unsqueeze(1).unsqueeze(-1)
num\ heads = 4
model = MultiHeadAttention(8, num_heads)
y = model(x, x, x, mask)
assert len(y.shape) == len(x.shape)
for dim_x, dim_y in zip(x.shape, y.shape):
   assert dim_x == dim_y
```

# 2.2 Positional Encoding Module

Positional Encoding is a critical component in the Transformer architecture, designed to provide information about the positions of elements in a sequence to a model that inherently lacks sequential information. Transformers use self-attention mechanisms that do not inherently understand the order or position of tokens in the input. Positional Encoding is introduced to address this limitation and allow the model to consider the order of elements within the input sequence.

It addresses the challenge of modeling sequences with self-attention mechanisms that do not inherently understand the order of elements. By adding Positional Encoding to the input embeddings, the model can differentiate between tokens based on their positions and capture sequential information effectively.

#### 1. Positional Encoding Function:

- Positional Encoding is typically represented as a fixed-size vector that is added element-wise to the input embeddings. This vector is determined by a mathematical function.
- The most common approach is to use a combination of sine and cosine functions with different frequencies and phases to create a unique encoding for each position.
- $\circ~$  For each position pos and dimension i of the Positional Encoding vector, PE(pos,2i) is given by

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

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

 $d_{model}$  is the dimension of the model's input embeddings.

#### 2. Adding Positional Encoding:

• The Positional Encoding vector is added element-wise to the input embeddings. This combination of the original word embeddings and the Positional Encoding allows the model to distinguish between tokens based on their positions.

#### 2.2.1 Let's Try to work an example!

Assume the sentence coming into the encoding is, "This is an Example"; The Positional Encoding layer is initialized with the following parameters;

- $k = 0 \le k < L = 3$
- n = 100
- d = 6
- $I = 0 \le i < d/2 = 2$
- Based on the values above, for the text given, Find the Position Encoding values below or create on on your own and upload to this section! DONT CODE IT

		Positional Encodings d = 6 & n = 100							
Sequence	Index (k)	i=0	i=0	i=1	i=1	i=2	i=2	i=3	i=3
This	0	$P_{00}=0$	$P_{01} = 1$	$P_{02} = 0$	$P_{03} = 1$	$P_{04} = 0$	$P_{05} = 1$		
is	1	$P_{10}$ =0.841	$P_{11} = 0.540$	$P_{12} = 0.002$	$P_{13} = 0.999$	$P_{14} = 4.641$	$P_{15} = 0.999$		
an	2	$P_{20} = 0.909$	$P_{21}$ = -0.416	$P_{22} = 0.004$	$P_{23} = 0.999$	$P_{24} = 9.283$	$P_{25} = 0.999$		
Example	3	$P_{30} = 0.141$	$P_{31}$ = -0.989	$P_{32} = 0.006$	$P_{33}$ = 0.999	$P_{34} = 1.392$	$P_{35}$ = 0.999		

#### ▼ 2.2.1 Let's Code!

Now try to use the same approach to implement the Positional Encoding part.

For full credit do not use for loops;

Make use of packages like

- torch.arange(): <u>https://pytorch.org/docs/stable/generated/torch.arange.html</u>
- torch.stack(): <a href="https://pytorch.org/docs/stable/generated/torch.stack.html">https://pytorch.org/docs/stable/generated/torch.stack.html</a>

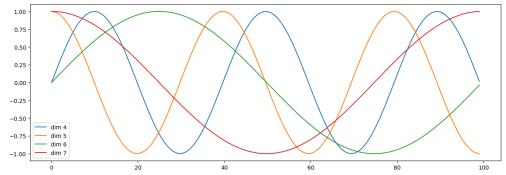
```
class PositionalEncoding(nn.Module):
   A module that adds positional encoding to each of the token's features.
   So that the Transformer is position aware.
   def
        _init__(self, input_dim: int, max_len: int=10000):
       - input_dim: Input dimension about the features for each token
       - max_len: The maximum sequence length
       super(PositionalEncoding, self).__init__()
       self.input_dim = input_dim
       self.max_len = max_len
   def forward(self, x):
       Compute the positional encoding and add it to x.
       - x: Tensor of the shape BxLxC, where B is the batch size, L is the sequence length,
        and {\sf C} is the channel dimension
       Return:
       - x: Tensor of the shape BxLxC, with the positional encoding added to the input
       seq_len = x.shape[1]
       input_dim = x.shape[2]
       pe = None
       # TODO: Compute the positional encoding
       # Check Section 3.5 for the definition (https://arxiv.org/pdf/1706.03762.pdf)
       # It's a bit messy, but the definition is provided for your here for your #
       # convenience (in LaTex).
       # PE_{(pos,2i)} = sin(pos / 10000^{2i}\dmodel})
                                                                       #
       \# PE_{(pos,2i+1)} = cos(pos / 10000^{2i}\dmodel})
                                                                       #
                                                                       #
       # You should replace 10000 with max len here.
       # Calculate position indices and division terms
       position = torch.arange(seq_len).unsqueeze(1).float().to(x.device)
       div_term = torch.exp(torch.arange(0, input_dim, 2).float() * (-math.log(self.max_len) / input_dim)).to(x.device)
       # Initialize the positional encoding matrix
       pe = torch.zeros(seq_len, input_dim).to(x.device)
       # Compute sine and cosine for positional encoding in one step
       pe[:, 0::2] = torch.sin(position * div_term)
       pe[:, 1::2] = torch.cos(position * div_term)
       # Replicate for batch size
       pe = pe.expand(x.size(0), -1, -1)
       END OF YOUR CODE
       x = x + pe.to(x.device)
       return x
# Sanity check - I
x = torch.randn(1, 100, 20)
pe = PositionalEncoding(20)
y = pe(x)
assert len(x.shape) == len(y.shape)
for dim_x, dim_y in zip(x.shape, y.shape):
   assert dim_x == dim_y
```

```
# Sanity Check - II
x = torch.randn(1, 100, 6)
n = 100
pe = PositionalEncoding(d,n)
y = pe(x)
y -= x
print(y[:,:4,:])
     tensor([[[ 0.0000, 1.0000, 0.0000, 1.0000,
                                                     0.0000,
                                           0.9769,
                                                     0.0464,
                                                              0.9989],
               0.8415, 0.5403, 0.2138,
               0.9093, -0.4161,
                                  0.4177,
                                           0.9086,
                                                     0.0927,
                                                              0.9957],
              [ 0.1411, -0.9900, 0.6023,
                                           0.7983,
                                                     0.1388,
                                                              0.9903]]])
tensor([[[ 0.0000, 1.0000, 0.0000,
                                   1.0000,
                                            0.0000.
          0.8415, 0.5403, 0.2138, 0.9769,
                                            0.0464,
          0.9093, -0.4161, 0.4177,
                                   0.9086,
                                            0.0927,
                                                    0.9957],
          0.1411, -0.9900,
                           0.6023,
                                    0.7983,
                                            0.1388,
                                                     0.9903]]])
```

```
# Sanity check - III
%matplotlib inline
import matplotlib.pyplot as plt
import numpy as np

plt.figure(figsize=(15, 5))
pe = PositionalEncoding(20)
y = pe.forward((torch.zeros(1, 100, 20)))
plt.plot(np.arange(100), y[0, :, 4:8].data.numpy())
plt.legend(["dim %d"%p for p in [4,5,6,7]])
```

<matplotlib.legend.Legend at 0x7987c064d210>



# 2.3 FeedForward Module

The FeedForward Layer in the Transformer architecture is a position-wise neural network layer designed to process the context-aware representations generated by the self-attention mechanism. It consists of two linear transformations followed by a non-linear activation function, typically ReLU. The FeedForward Layer is applied independently to each position in the sequence, allowing the model to capture different patterns at different positions. This position-wise independence, combined with non-linearity, helps the model learn complex relationships within the data and plays a crucial role in the Transformer's ability to process and understand sequential data effectively, making it a fundamental component for various sequence-to-sequence tasks.

Mathematically, if X represents the input sequence (a sequence of embeddings), FFN(X) is the output of the FeedForward Layer, and  $W_1$ ,  $W_2$ ,  $b_1$ , and  $b_2$  represent learned weight matrices and bias terms, the operation can be expressed as,

$$FFN(X) = ReLU(X. W_i + b_i). W_2 + b_2.$$

```
class FeedForwardNetwork(nn.Module):
  A simple feedforward network. Essentially, it is a two-layer fully-connected
  neural network.
  def
     _init__(self, input_dim, ff_dim):
     Inputs:
     - input_dim: Input dimension
     - ff dim: Hidden dimension
     super(FeedForwardNetwork, self).__init__()
     # TODO: Define the two linear layers and a non-linear one.
     self.linear1 = nn.Linear(input_dim, ff_dim)
     self.relu = nn.ReLU()
     self.linear2 = nn.Linear(ff_dim, input_dim)
     END OF YOUR CODE
     def forward(self, x: torch.Tensor):
     Input:
     - x: Tensor of the shape BxLxC, where B is the batch size, L is the sequence length,
     and C is the channel dimension
     Return:
     - y: Tensor of the shape BxLxC
     y = None
     # TODO: Process the input.
     y = self.linear2(self.relu(self.linear1(x)))
     END OF YOUR CODE
     return y
# Sanity Check
x = torch.randn((2, 10, 8))
ff dim = 4
model = FeedForwardNetwork(8, ff_dim)
y = model(x)
assert len(x.shape) == len(y.shape)
for dim_x, dim_y in zip(x.shape, y.shape):
  assert dim_x == dim_y
print(y.shape)
  torch.Size([2, 10, 8])
```

#### 2.4 Encoder Module

The Encoder module in a Transformer is responsible for processing the input sequence, typically used for tasks like language understanding and representation learning. It consists of multiple identical layers, each containing two main components: the Multi-Head Self-Attention mechanism and the Position-wise FeedForward Layer.

- In each layer, the input sequence is first passed through the Multi-Head Self-Attention mechanism, which computes weighted representations for each element in the sequence, capturing contextual information. The attention output is then passed through the Position-wise FeedForward Layer, introducing non-linearity and allowing the model to capture different patterns at each position.
- This process is repeated for each layer in the encoder stack, enabling the model to capture hierarchical features and within the input sequence effectively. The final encoder output represents a rich contextualized representation of the input sequence, which can be used for various downstream tasks, including translation, text generation, and sentiment analysis.

#### 2.4.1 Encoder Cell

```
class TransformerEncoderCell(nn.Module):
   A single cell (unit) for the Transformer encoder.
   def __init__(self, input_dim: int, num_heads: int, ff_dim: int, dropout: float):
      Inputs:
      - input_dim: Input dimension for each token in a sequence
      - num_heads: Number of attention heads in a multi-head attention module
      - ff dim: The hidden dimension for a feedforward network
      - dropout: Dropout ratio for the output of the multi-head attention and feedforward
       modules.
      super(TransformerEncoderCell, self).__init__()
      # TODO: A single Transformer encoder cell consists of
      # 1. A multi-head attention module
      # 2. Followed by dropout
      # 3. Followed by layer norm (check nn.LayerNorm)
      # https://pytorch.org/docs/stable/generated/torch.nn.LayerNorm.html#torch.nn.LayerNorm
      # At the same time, it also has
      # 1. A feedforward network
      # 2. Followed by dropout
      # 3. Followed by layer norm
      # Initialize multi-head attention and feedforward network
      self.multihead_attn = MultiHeadAttention(input_dim, num_heads)
      self.ff_network = FeedForwardNetwork(input_dim, ff_dim)
      # Initialize dropout and layer normalization for both parts of the cell
      self.attn_dropout = nn.Dropout(dropout)
      self.ffn_dropout = nn.Dropout(dropout)
      self.attn_norm = nn.LayerNorm(input_dim)
      self.ffn_norm = nn.LayerNorm(input_dim)
      # Placeholder for attention output (for visualization purposes)
      self.attention_output = None
      END OF YOUR CODE
      self.attention = None
   def forward(self, x: torch.Tensor, mask: torch.Tensor=None):
      Inputs:
      - x: Tensor of the shape BxLxC, where B is the batch size, L is the sequence length,
        and {\sf C} is the channel dimension
       mask: Tensor for multi-head attention
      y = None
      # TODO: Get the output of the multi-head attention part (with dropout
      # and layer norm), which is used as input to the feedforward network (
                                                                  #
      # again, followed by dropout and layer norm).
                                                                  #
                                                                  #
      # Don't forget the residual connections for both parts. Append the
      # 1st Normalized Output before feed_forward to self.attention(Useful in
                                                                  #
      # visualizing)
      # Apply multi-head attention
      attn_output = self.multihead_attn(x, x, x, mask)
      # Apply dropout and layer norm with residual connection
      attn_output = self.attn_dropout(attn_output)
      attn_output = self.attn_norm(attn_output + x)
      # Store the attention output (optional, for visualization)
      self.attention_output = attn_output
      # Apply feedforward network
      ffn_output = self.ff_network(attn_output)
      # Apply dropout and layer norm with residual connection
      ffn_output = self.ffn_dropout(ffn_output)
      y = self.ffn_norm(ffn_output + attn_output)
      END OF YOUR CODE
```

```
# Sanity Check

x = torch.randn((2, 10, 8))
mask = torch.randn((2, 10)) > 0.5
mask = mask.unsqueeze(1).unsqueeze(-1)
num_heads = 4
model = TransformerEncoderCell(8, num_heads, 32, 0.1)
y = model(x, mask)
assert len(x.shape) == len(y.shape)
for dim_x, dim_y in zip(x.shape, y.shape):
    assert dim_x == dim_y
```

▼ 2.4.2 Building an Encoder Module

torch.Size([2, 10, 8])

return y

print(y.shape)

```
class TransformerEncoder(nn.Module):
   A full encoder consisting of a set of TransformerEncoderCell.
  def __l
      _init__(self, input_dim: int, num_heads: int, ff_dim: int, num_cells: int, dropout: float=0.1):
      Inputs:
      - input_dim: Input dimension for each token in a sequence
      - num_heads: Number of attention heads in a multi-head attention module
      - ff dim: The hidden dimension for a feedforward network
      num_cells: Number of TransformerEncoderCells
      - dropout: Dropout ratio for the output of the multi-head attention and feedforward
      super(TransformerEncoder, self).__init__()
      self.norm = None # LaverNorm laver
      self.cells = None # TransformerEncoderCells
      # TODO: Construct a nn.ModuleList to store a stack of
      # TranformerEncoderCells. Check the documentation here of how to use it
      # https://pytorch.org/docs/stable/generated/torch.nn.ModuleList.html#torch.nn.ModuleList
      # At the same time, define a layer normalization layer to process the
      # output of the entire encoder.
      # Initialize a list of Transformer encoder cells
      self.encoder\_cells = nn. ModuleList([TransformerEncoderCell(input\_dim, num\_heads, ff\_dim, dropout) \ for \ \_in \ range(num\_cells)])
      # Initialize a layer normalization layer
      self.final norm = nn.LaverNorm(input dim)
      END OF YOUR CODE
      def forward(self, x: torch.Tensor, mask: torch.Tensor=None):
      Inputs:
      - x: Tensor of the shape BxLxC, where B is the batch size, L is the sequence length,
       and C is the channel dimension
      - mask: Tensor for multi-head attention
      Return:
      - y: Tensor of the shape of BxLxC, which is the normalized output of the encoder
      y = None
      # TODO: Feed x into the stack of TransformerEncoderCells and then
      # normalize the output with layer norm.
      # Pass the input through each encoder cell
      for encoder_cell in self.encoder_cells:
         x = encoder_cell(x, mask)
      # Normalize the output of the final encoder cell
      y = self.final_norm(x)
      END OF YOUR CODE
      return y
```

# 2.5 Transformer Classifier

Now, lets put this all the above describled modules together to make out classifier

```
A Transformer-based text classifier.
def
   __init__(self,
      vocab_size: int, embed_dim: int, num_heads: int, trx_ff_dim: int,
      num_trx_cells: int, num_class: int, dropout: float=0.1, pad_token: int=0
   ):
   .....
   Inputs:
   - vocab_size: Vocabulary size, indicating how many tokens we have in total.
   - embed_dim: The dimension of word embeddings
   - num_heads: Number of attention heads in a multi-head attention module
   - trx_ff_dim: The hidden dimension for a feedforward network
   - num_trx_cells: Number of TransformerEncoderCells
   - dropout: Dropout ratio
   - pad_token: The index of the padding token.
   super(TransformerClassifier, self).__init__()
   self.embed_dim = embed_dim
   # word embedding layer
   self.embedding = nn.Embedding(vocab_size, embed_dim, padding_idx=pad_token)
   # TODO: Define a module for positional encoding, Transformer encoder, and #
   # a output layer
   # Positional encoding
   self.pos_encoding = PositionalEncoding(embed_dim)
   # Transformer encoder
   self.transformer_encoder = TransformerEncoder(embed_dim, num_heads, trx_ff_dim, num_trx_cells, dropout)
   # Output linear layer
   self.classification_head = nn.Linear(embed_dim, num_class)
   END OF YOUR CODE
   def forward(self, text, mask=None):
   Inputs:
   - text: Tensor with the shape of BxLxC.
   - mask: Tensor for multi-head attention
   Return:
   - logits: Tensor with the shape of BxK, where K is the number of classes
   # word embeddings, note we multiple the embeddings by a factor
   embedded = self.embedding(text) * math.sqrt(self.embed_dim)
   logits = None
   # TODO: Apply positional embedding to the input, which is then fed into
   # the encoder. Average pooling is applied then to all the features of all #
   \# tokens. Finally, the logits are computed based on the pooled features. \#
   # Apply positional encoding
   encoded_tokens = self.pos_encoding(embedded)
   transformer_output = self.transformer_encoder(encoded_tokens, mask)
   sequence_avg = transformer_output.mean(dim=1)
   logits = self.classification_head(sequence_avg)
   END OF YOUR CODE
```

return logits

class TransformerClassifier(nn.Module):

```
# Sanity Check
vocab size = 10
embed_dim = 16
num\ heads = 4
trx_ff_dim = 16
num_trx_cells = 2
num_class = 3
x = torch.arange(vocab_size).view(1, -1)
x = torch.cat((x, x), dim=0)
mask = (x != 0).unsqueeze(-2).unsqueeze(1)
model = TransformerClassifier(vocab_size, embed_dim, num_heads, trx_ff_dim, num_trx_cells, num_class)
print('x: {}, mask: {}'.format(x.shape, mask.shape))
y = model(x, mask)
assert len(y.shape) == 2 and y.shape[0] == x.shape[0] and y.shape[1] == num_class
print(y.shape)
    x: torch.Size([2, 10]), mask: torch.Size([2, 1, 1, 10])
    torch.Size([2, 3])
```

# 2.6 Deciding HyperParameters & Training

```
assert torch.cuda.is_available()
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
# Hyperparameters
epochs = 3 \# epoch
lr = 0.0005 # learning rate
batch_size = 64 # batch size for training
train_iter = AG_NEWS(split='train')
num_class = len(set([label for (label, text) in train_iter]))
vocab_size = len(vocab)
emsize = 64
num_heads = 4
num_trx_cells = 2
gradient_norm_clip = 1
# Define a Transformer-based text classifier and a loss function.
# Define the criterion for classification
loss_func = nn.CrossEntropyLoss() # <-- be careful with the name</pre>
# Define the model
model = TransformerClassifier( # <-- be careful with the name</pre>
   vocab_size=vocab_size,
   embed_dim=emsize,
   num_heads=num_heads,
   trx_ff_dim=ff_dim,
   num_trx_cells=num_trx_cells,
   num_class=num_class,
   pad_token=PAD_TOKEN)
END OF YOUR CODE
model = model.to(device)
optimizer = torch.optim.AdamW(model.parameters(), lr=lr)
scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(optimizer, epochs, 1e-8)
total_accu = None
# You should be able to get a validation accuracy around 89%
for epoch in range(1, epochs + 1):
   epoch_start_time = time.time()
   train(model, train_dataloader, loss_func, device, gradient_norm_clip, optimizer)
   accu_val = evaluate(model, valid_dataloader, loss_func, device)
   if total_accu is not None and total_accu > accu_val:
      scheduler.step()
   else:
       total_accu = accu_val
   print('-' * 59)
   print('| end of epoch {:3d} | time: {:5.2f}s | '
        'valid accuracy {:8.3f} '.format(epoch,
                                    time.time() - epoch_start_time,
                                    accu_val))
   print('-' * 59)
```

	epoch epoch epoch	1   1   1	1000/	1782	batches batches batches		accuracy accuracy accuracy	0.643 0.794 0.830	
I	end of	epoch	1	time	: 23.16s		valid accu	ıracy	0.848
	epoch epoch epoch	2   2   2	1000/	1782	batches batches batches		accuracy accuracy accuracy	0.868 0.879 0.883	
1	end of	epoch	2	time	: 23.42s		valid accu	ıracy	0.874
	epoch epoch epoch	3   3   3	1000/	1782	batches batches batches	   	accuracy accuracy accuracy	0.902 0.903 0.903	
	end of	epoch	3	time	: 23.38s		valid accu	ıracy	0.886

# → Section 3: Transformers - 'Attention is All you need' : Machine Translation



The decoder module in a Seq-Seq Transformer model is responsible for generating the output sequence based on the information gathered by the encoder and previous tokens in an autoregressive manner. It utilizes self-attention mechanisms with masking to enforce causality and multi-head attention to capture dependencies between tokens in the output. Self-attention calculates attention weights for each position in the sequence, and multi-head attention aggregates the results from multiple attention heads, enhancing the model's representational power. Cross-attention is also employed to allow the decoder to focus on relevant parts of the encoder's output.

Additionally, position-wise feed-forward networks further process the information by applying linear transformations and non-linear activation functions to each position independently. Throughout the decoder, layer normalization and residual connections are utilized to enhance training stability. These components are typically stacked in multiple layers to enable the model to learn complex relationships and generate coherent output sequences.

#### ▼ 3.1.1 Decoder Cell

```
A single cell (unit) of the Transformer decoder.
def __l
    _init__(self, input_dim: int, num_heads: int, ff_dim: int, dropout: float=0.1):
   Inputs:
   - input_dim: Input dimension for each token in a sequence
   - num_heads: Number of attention heads in a multi-head attention module
   - ff dim: The hidden dimension for a feedforward network
   - dropout: Dropout ratio for the output of the multi-head attention and feedforward
     modules.
   super(TransformerDecoderCell, self).__init__()
   # TODO: Similar to the TransformerEncoderCell. define two
                                                                #
   # MultiHeadAttention modules. One for processing the tokens on the
                                                                 #
                                                                 #
   # decoder side. The other for getting the attention across the encoder.
   # and the decoder. Also define a feedforward network. Don't forget the
                                                                 #
   # Dropout and Layer Norm layers.
   # Self-attention mechanism
   self.decoder_self_attention = MultiHeadAttention(input_dim, num_heads)
   self.self_attn_dropout = nn.Dropout(dropout)
   self.self_attn_norm = nn.LayerNorm(input_dim)
   # Cross-attention mechanism
   self.decoder_cross_attention = MultiHeadAttention(input_dim, num_heads)
   self.cross_attn_dropout = nn.Dropout(dropout)
   self.cross_attn_norm = nn.LayerNorm(input_dim)
   # Feedforward network
   self.decoder_ffn = FeedForwardNetwork(input_dim, ff_dim)
   self.ffn_dropout = nn.Dropout(dropout)
   self.ffn_norm = nn.LayerNorm(input_dim)
   END OF YOUR CODE
   def forward(self, x: torch.Tensor, encoder_output: torch.Tensor, src_mask=None, tgt_mask=None):
   Inputs:
   - x: Tensor of BxLdxC, word embeddings on the decoder side
   - encoder_output: Tensor of BxLexC, word embeddings on the encoder side
   - src_mask: Tensor, masks of the tokens on the encoder side
   - tgt_mask: Tensor, masks of the tokens on the decoder side
   Return:
   - y: Tensor of BxLdxC. Attended features for all tokens on the decoder side.
   .....
   y = None
   # TODO: Compute the self-attended features for the tokens on the decoder #
   # side. Then compute the corss-attended features for the tokens on the
   # decoder side to the encoded features, which are finally feed into the
   # feedforward network
   # Self-attention on decoder input
   decoder_self_attended = self.decoder_self_attention(x, x, x, tgt_mask)
   decoder_self_attended = self.self_attn_dropout(decoder_self_attended)
   decoder_self_attended = self.self_attn_norm(decoder_self_attended + x)
   # Cross-attention with encoder output
   decoder_cross_attended = self.decoder_cross_attention(decoder_self_attended, encoder_output, encoder_output, src_mask)
   decoder_cross_attended = self.cross_attn_dropout(decoder_cross_attended)
   decoder_cross_attended = self.cross_attn_norm(decoder_cross_attended + decoder_self_attended)
   # Pass through the feedforward network
   ffn_output = self.decoder_ffn(decoder_cross_attended)
   ffn_output = self.ffn_dropout(ffn_output)
   y = self.ffn_norm(ffn_output + decoder_cross_attended)
   END OF YOUR CODE
```

class TransformerDecoderCell(nn.Module):

```
dec_feats = torch.randn((3, 10, 16))
dec_mask = torch.randn((3, 1, 10, 10)) > 0.5

enc_feats = torch.randn((3, 12, 16))
enc_mask = torch.randn((3, 1, 1, 12)) > 0.5

model = TransformerDecoderCell(16, 2, 32, 0.1)
z = model(dec_feats, enc_feats, enc_mask, dec_mask)
assert len(z.shape) == len(dec_feats.shape)
for dim_z, dim_x in zip(z.shape, dec_feats.shape):
    assert dim_z == dim_x
print(z.shape)

    torch.Size([3, 10, 16])
```

# ▼ 3.1.2 Building the Decoder Module

# Sanity Check

```
class TransformerDecoder(nn.Module):
    A TransformerDecoder is a stack of multiple TransformerDecoderCells and a Layer Norm.
# Sanity Check
dec_feats = torch.randn((3, 10, 16))
dec_mask = torch.randn((3, 1, 10, 10)) > 0.5
enc_feats = torch.randn((3, 1, 1, 12)) > 0.5

model = TransformerDecoder(16, 2, 32, 2, 0.1)
z = model(dec_feats, enc_feats, enc_mask, dec_mask)
assert len(z.shape) == len(dec_feats.shape)
for dim_z, dim_x in zip(z.shape, dec_feats.shape):
    assert dim_z == dim_x
print(z.shape)
    torch.Size([3, 10, 16])
```

# 3.2 Transformer Based Seq-to-Seq Model

```
# Croate a stack of TransformerDecoderColle
class Seq2SeqTransformer(nn.Module):
   Transformer-based sequence-to-sequence model.
   def __init__(self,
          num_encoder_layers: int, num_decoder_layers: int, embed_dim: int,
          num_heads: int, src_vocab_size: int, tgt_vocab_size: int,
          trx_ff_dim: int = 512, dropout: float = 0.1, pad_token: int=0
      ):
      .....
      Inputs:
      - num_encoder_layers: How many TransformerEncoderCell in stack
      - num_decoder_layers: How many TransformerDecoderCell in stack
      - embed_dim: Word embeddings dimension
      - num_heads: Number of attention heads
      - src_vocab_size: Number of tokens in the source language vocabulary
      tgt_vocab_size: Number of tokens in the target language vocabulary
      - trx_ff_dim: Hidden dimension in the feedforward network
      - dropout: Dropout ratio
      super(Seq2SeqTransformer, self).__init__()
      self.embed_dim = embed_dim
      # Word embeddings for both the source and target languages
      self.src_token_embed = nn.Embedding(src_vocab_size, embed_dim, padding_idx=pad_token)
      self.tgt_token_embed = nn.Embedding(tgt_vocab_size, embed_dim, padding_idx=pad_token)
      # TODO: Define the positional encoding, encoder, decoder, and the output #
      # layer. Think of how many classes are in the output layer.
      # Positional Encoding
      self.positional_encoding = PositionalEncoding(embed_dim)
      # Transformer Encoder and Decoder
      self.transformer_encoder = TransformerEncoder(embed_dim, num_heads, trx_ff_dim, num_encoder_layers, dropout)
      self.transformer_decoder = TransformerDecoder(embed_dim, num_heads, trx_ff_dim, num_decoder_layers, dropout)
      # Output linear layer
      self.output_layer = nn.Linear(embed_dim, tgt_vocab_size)
```