

TRANSFORMER TRANSLATION

SG-1

TOKENIZATION

Tokenization involves converting the text supplied to the model into smaller parts.

For ex: I spent \$100 in the U.S.A.

The tokens of the above sentence would be
['i', 'spent', '\$', '100', 'in', 'the', 'U.S.A']

Here the tokenisation was done by converting entire words into tokens. Sometimes this isn't the best form of tokenisation and subword tokenisation would be better by making parts of words that repeat frequently to instead be the tokens rather than every new word being a new token. Helps deal with out of vocab words.

POSITIONAL ENCODING

It encodes information about the position of the token in the sentence by a bunch of sines and cosines.

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

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

It is a function of the position of the word and dimension i in the N dimensional position embedding.

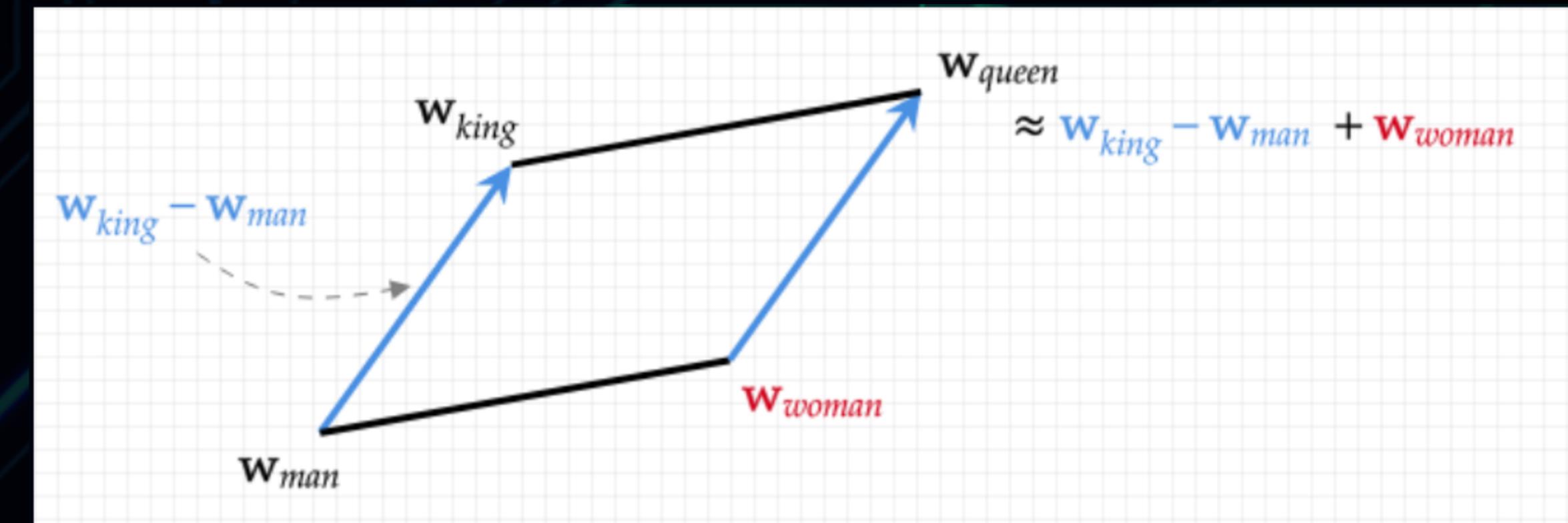
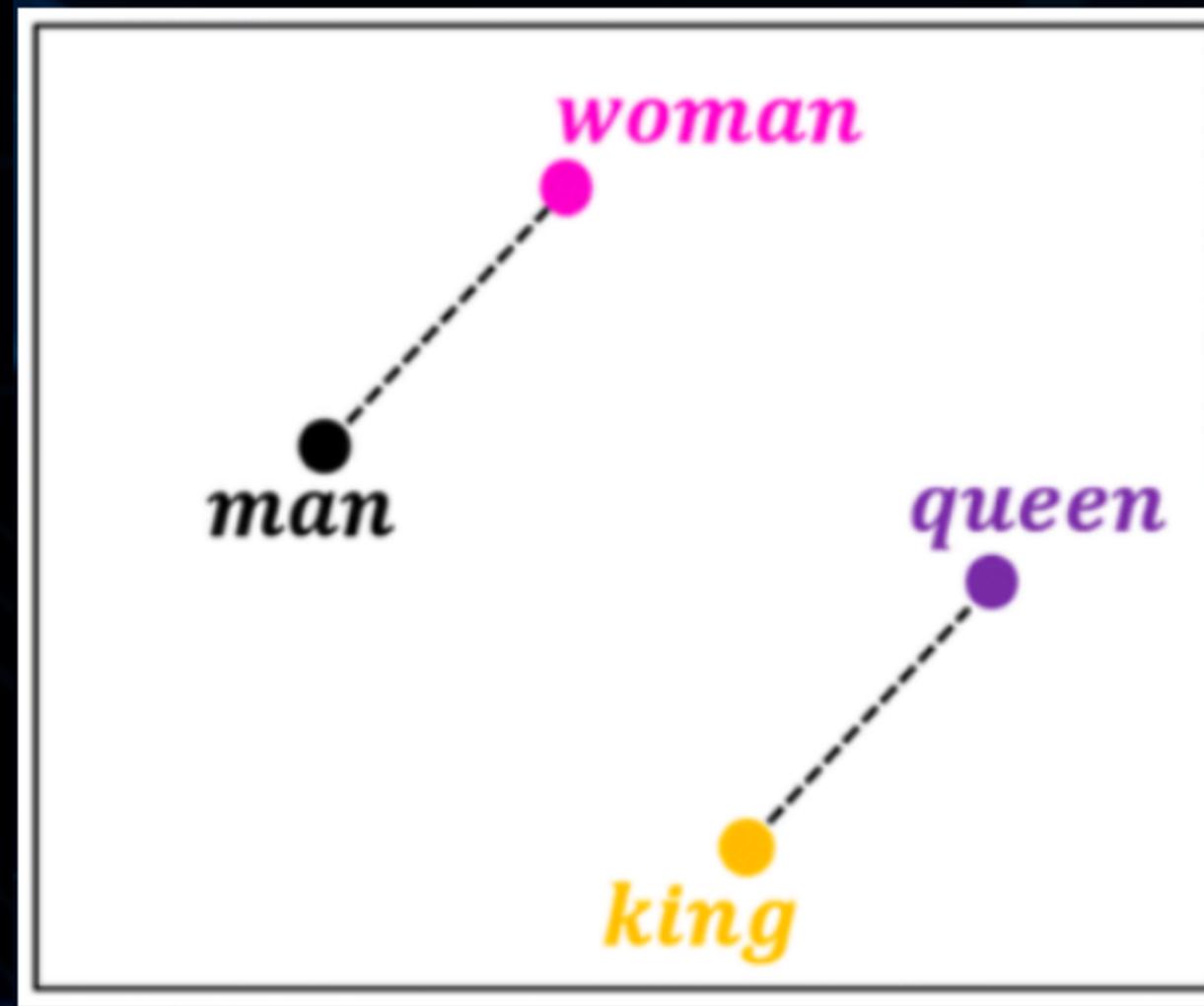
EMBEDDINGS

Adds meaning to the tokens by encoding data from position/location of the words in context.

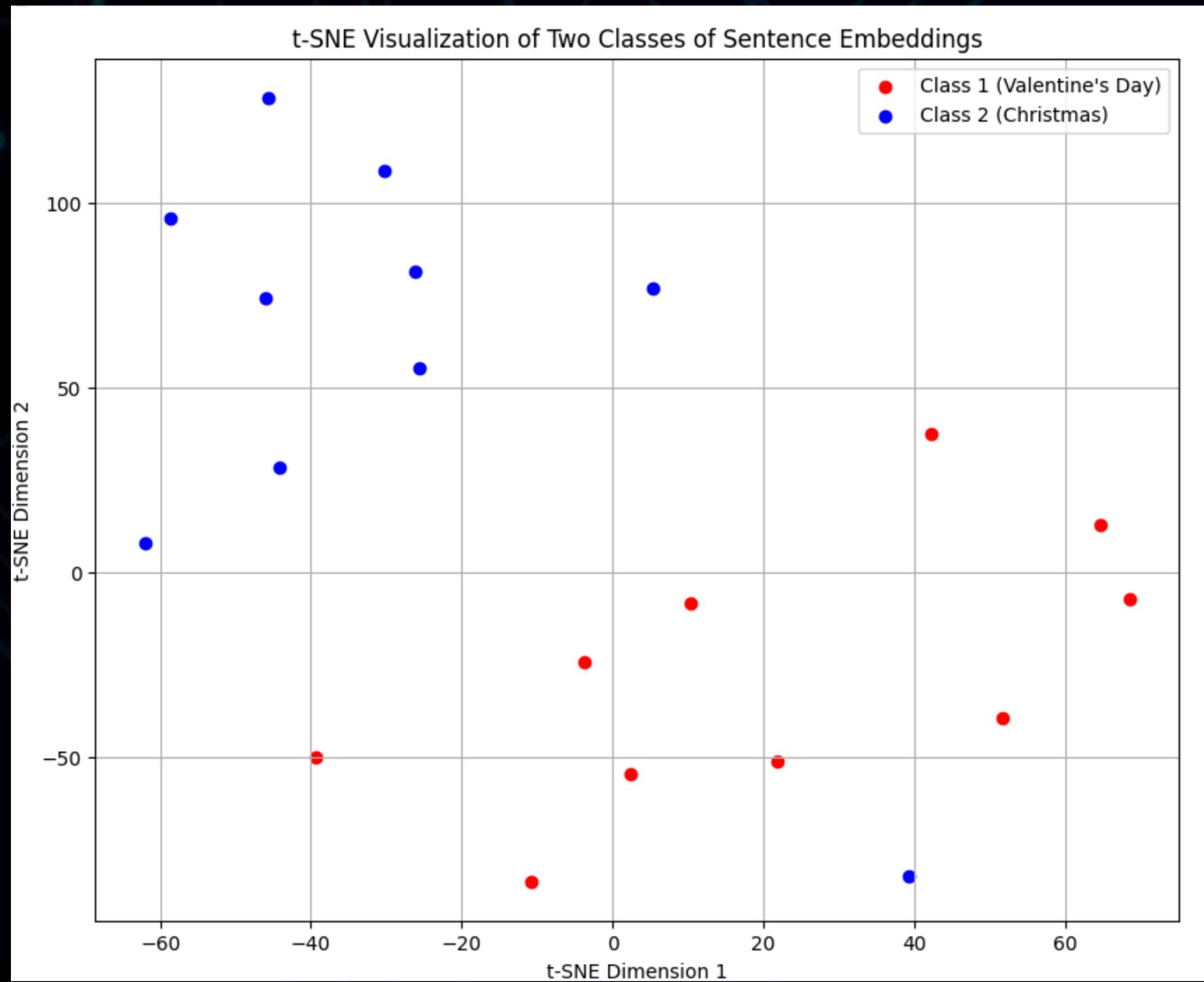
Embeddings are what give meaning to the tokens
two different tokens can have similar embeddings

For ex: King and Queen would have similar embeddings
and if we give ‘King’ - ‘man’ + ‘woman’, it would return ‘Queen’

WORD EMBEDDINGS VISUALISED

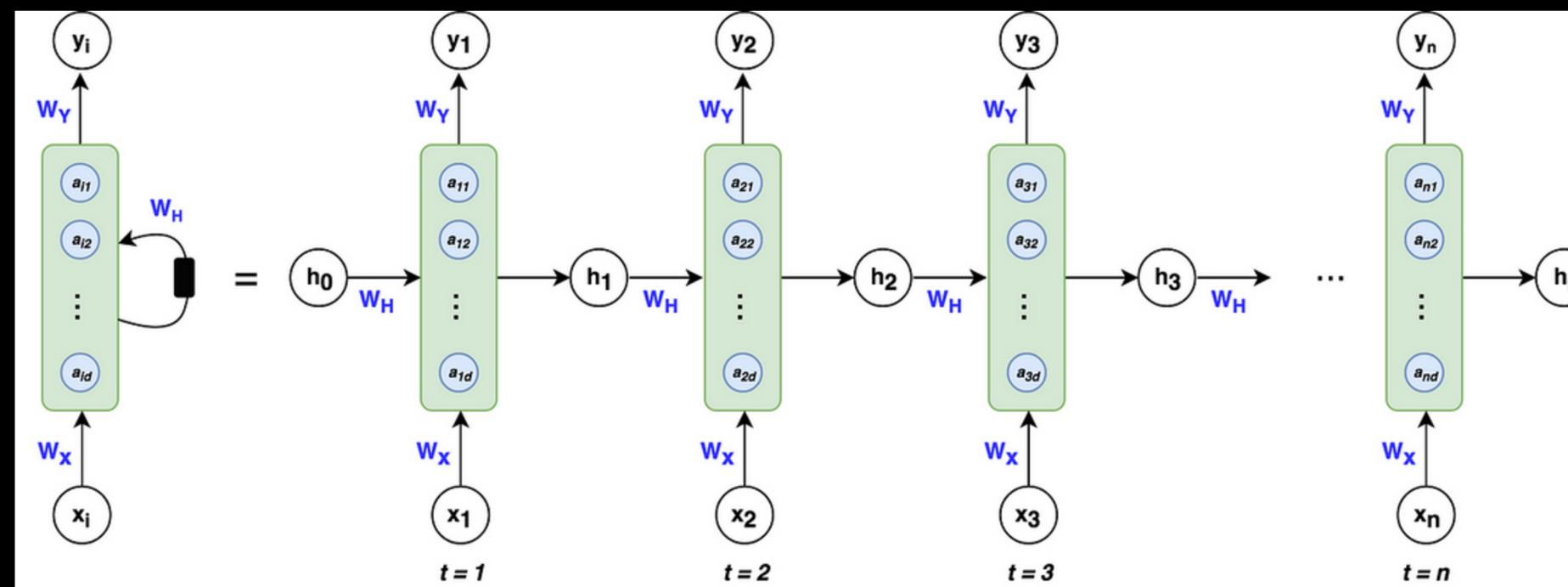


CLUSTERING



RNN

To handle sequential data in NLPs,
we use RNNs.



- Takes in account the order of words
- Flexible no. of inputs rather than fixed no. of inputs
- Preserves contextual meaning

Issues: Exploding/vanishing gradients because of which it is not suitable for long sequences.

LSTM

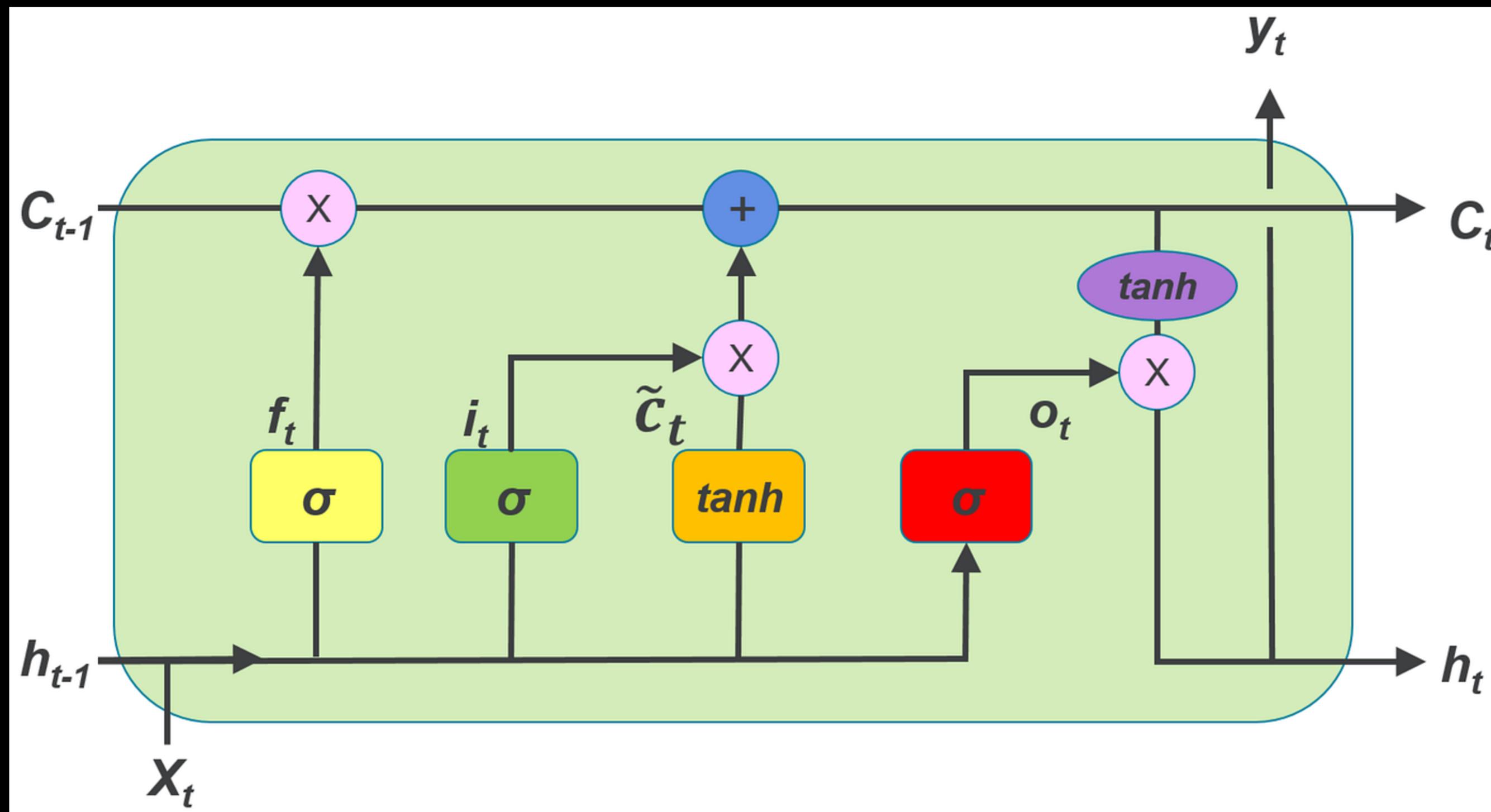
LSTMs are designed to avoid the issue of exploding/ vanishing gradients.

Instead of having the same loop for both past and recent events, it has two path, one to update the long term memory and one to update the short term memory.

The LSTM structure consists of 3 gates:

- Forget Gate: Decides what percentage of the existing long term memory to be remembered
- Input Gate: Decides what percentage of the new information to be stored
- Output Gate: Controls the output flow of the memory content to the next layer

STRUCTURE OF LSTM



ATTENTION

WHY ATTENTION

In the case of transformers, we want words to not only represent their individual meaning, but also soak in the context of the sentence. It is the job of an attention block to calculate what it needs to add to the generic embedding, as a function of its context, to move it to one of those specific directions.

ATTENTION BLOCK

- Query
- Key
- Value

MLP BLOCK

- linear
- ReLU
- linear

ATTENTION

QUERY

Query vector is obtained by multiplying W_q matrix with the embedding vector for each token. It can be thought of as each token asking a question

KEY

Key vector is obtained by multiplying W_k matrix with the embedding vector for each token. It can be thought of as each query being answered.

We take the dot product of the key and query vector to find out how closely each query-key pair aligns with each other. We then apply softmax, and we can think of the values we get as the weight of how relevant one word is to the other. We may apply masking as well.

ATTENTION

VALUE

This W_v matrix is multiplied by every one of those embeddings to produce a sequence of value vectors. Finally, these value vectors are scaled by the weights we obtained previously. It tells us what specific adjustment to make so that the context is also absorbed in it.

MULTI HEADED ATTENTION

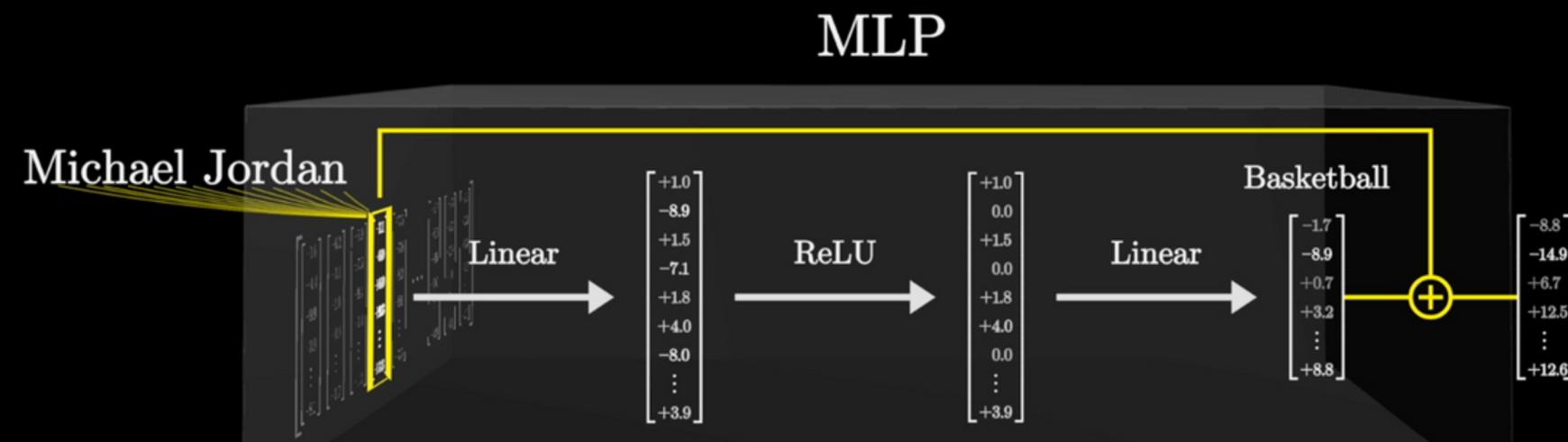
it can be thought of as repeating the self attention blocks in parallel. the results obtained from each of the self attention heads are then summed and added to the embedding vector.

ATTENTION

MULTI LAYER PERCEPTRON

it is responsible for storing facts.

“Michael Jordan plays Basketball”



FIRST LINEAR LAYER

multiplication with a weight matrix,
which can be thought of as each row
asking a question to the entries of
the embedding vector

activation of neuron tells us how
much a vector aligns with a
particular direction

CODE IMPLEMENTATION

FUNCTIONS FOR CREATING VOCABULARY AND TOKENIZATION :

```
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

SRC_LANGUAGE = 'de'
TGT_LANGUAGE = 'en'

# Place-holders
token_transform = {}
vocab_transform = {}

# Create source and target language tokenizer. Make sure to install the dependencies.
# pip install -U spacy
! python -m spacy download en_core_web_sm
! python -m spacy download de_core_news_sm
token_transform[SRC_LANGUAGE] = get_tokenizer('spacy', language='de_core_news_sm')
token_transform[TGT_LANGUAGE] = get_tokenizer('spacy', language='en_core_web_sm')

# helper function to yield list of tokens
def yield_tokens(data_iter: Iterable, language: str) -> List[str]:
    language_index = {SRC_LANGUAGE: 0, TGT_LANGUAGE: 1}

    for data_sample in data_iter:
        yield token_transform[language](data_sample[language_index[language]])

# Define special symbols and indices
UNK_IDX, PAD_IDX, BOS_IDX, EOS_IDX = 0, 1, 2, 3
# Make sure the tokens are in order of their indices to properly insert them in vocab
special_symbols = ['<unk>', '<pad>', '<bos>', '<eos>']
```

CREATING VOCABULARY:

```
for ln in [SRC_LANGUAGE, TGT_LANGUAGE]:  
    # Training data Iterator  
    train_iter = Multi30k(split='train', language_pair=(SRC_LANGUAGE, TGT_LANGUAGE))  
    # Create torchtext's Vocab object  
    vocab_transform[ln] = build_vocab_from_iterator(yield_tokens(train_iter, ln),  
                                                    min_freq=1,  
                                                    specials=special_symbols,  
                                                    special_first=True)  
  
    # Set UNK_IDX as the default index. This index is returned when the token is not found.  
    # If not set, it throws RuntimeError when the queried token is not found in the Vocabulary.  
    for ln in [SRC_LANGUAGE, TGT_LANGUAGE]:  
        vocab_transform[ln].set_default_index(UNK_IDX)
```

```
from torch import Tensor  
import torch  
import torch.nn as nn  
from torch.nn import Transformer  
import math  
DEVICE = torch.device('cuda' if torch.cuda.is_available() else 'cpu')  
  
# helper Module that adds positional encoding to the token embedding to introduce a notion of word order.  
class PositionalEncoding(nn.Module):  
    def __init__(self,  
                 emb_size: int,  
                 dropout: float,  
                 maxlen: int = 5000):  
        super(PositionalEncoding, self).__init__()  
        den = torch.exp(- torch.arange(0, emb_size, 2) * math.log(10000) / emb_size)  
        pos = torch.arange(0, maxlen).reshape(maxlen, 1)  
        pos_embedding = torch.zeros((maxlen, emb_size))  
        pos_embedding[:, 0::2] = torch.sin(pos * den)  
        pos_embedding[:, 1::2] = torch.cos(pos * den)  
        pos_embedding = pos_embedding.unsqueeze(-2)  
  
        self.dropout = nn.Dropout(dropout)  
        self.register_buffer('pos_embedding', pos_embedding)  
  
    def forward(self, token_embedding: Tensor):  
        return self.dropout(token_embedding + self.pos_embedding[:token_embedding.size(0), :])
```

: POSITIONAL ENCODING

CODE IMPLEMENTATION

CREATING EMBEDDINGS:

CORE TRANSFORMER ARCHITECTURE:

```
class TokenEmbedding(nn.Module):
    def __init__(self, vocab_size: int, emb_size):
        super(TokenEmbedding, self).__init__()
        self.embedding = nn.Embedding(vocab_size, emb_size)
        self.emb_size = emb_size

    def forward(self, tokens: Tensor):
        return self.embedding(tokens.long()) * math.sqrt(self.emb_size)

# Seq2Seq Network
class Seq2SeqTransformer(nn.Module):
    def __init__(self,
                 num_encoder_layers: int,
                 num_decoder_layers: int,
                 emb_size: int,
                 nhead: int,
                 src_vocab_size: int,
                 tgt_vocab_size: int,
                 dim_feedforward: int = 512,
                 dropout: float = 0.1):
        super(Seq2SeqTransformer, self).__init__()
        self.transformer = Transformer(d_model=emb_size,
                                      nhead=nhead,
                                      num_encoder_layers=num_encoder_layers,
                                      num_decoder_layers=num_decoder_layers,
                                      dim_feedforward=dim_feedforward,
                                      dropout=dropout)
        self.generator = nn.Linear(emb_size, tgt_vocab_size)
        self.src_tok_emb = TokenEmbedding(src_vocab_size, emb_size)
        self.tgt_tok_emb = TokenEmbedding(tgt_vocab_size, emb_size)
        self.positional_encoding = PositionalEncoding(
            emb_size, dropout=dropout)
```

CODE IMPLEMENTATION

CONTINUATION..

```
def forward(self,
            src: Tensor,
            trg: Tensor,
            src_mask: Tensor,
            tgt_mask: Tensor,
            src_padding_mask: Tensor,
            tgt_padding_mask: Tensor,
            memory_key_padding_mask: Tensor):
    src_emb = self.positional_encoding(self.src_tok_emb(src))
    tgt_emb = self.positional_encoding(self.tgt_tok_emb(trg))
    outs = self.transformer(src_emb, tgt_emb, src_mask, tgt_mask, None,
                           src_padding_mask, tgt_padding_mask, memory_key_padding_mask)
    return self.generator(outs)

def encode(self, src: Tensor, src_mask: Tensor):
    return self.transformer.encoder(self.positional_encoding(
        self.src_tok_emb(src)), src_mask)

def decode(self, tgt: Tensor, memory: Tensor, tgt_mask: Tensor):
    return self.transformer.decoder(self.positional_encoding(
        self.tgt_tok_emb(tgt)), memory,
                                   tgt_mask)
```

CODE IMPLEMENTATION

CREATING MASKS:

```
def generate_square_subsequent_mask(sz):
    mask = (torch.triu(torch.ones((sz, sz), device=DEVICE)) == 1).transpose(0, 1)
    mask = mask.float().masked_fill(mask == 0, float('-inf')).masked_fill(mask == 1, float(0.0))
    return mask

def create_mask(src, tgt):
    src_seq_len = src.shape[0]
    tgt_seq_len = tgt.shape[0]

    tgt_mask = generate_square_subsequent_mask(tgt_seq_len)
    src_mask = torch.zeros((src_seq_len, src_seq_len), device=DEVICE).type(torch.bool)

    src_padding_mask = (src == PAD_IDX).transpose(0, 1)
    tgt_padding_mask = (tgt == PAD_IDX).transpose(0, 1)
    return src_mask, tgt_mask, src_padding_mask, tgt_padding_mask
```

CODE IMPLEMENTATION

INITIALISING THE MODEL:

```
torch.manual_seed(10101100)

SRC_VOCAB_SIZE = len(vocab_transform[SRC_LANGUAGE])
TGT_VOCAB_SIZE = len(vocab_transform[TGT_LANGUAGE])
EMB_SIZE = 512
NHEAD = 8
FFN_HID_DIM = 512
BATCH_SIZE = 128
NUM_ENCODER_LAYERS = 3
NUM_DECODER_LAYERS = 3

transformer = Seq2SeqTransformer(NUM_ENCODER_LAYERS, NUM_DECODER_LAYERS, EMB_SIZE,
                                 NHEAD, SRC_VOCAB_SIZE, TGT_VOCAB_SIZE, FFN_HID_DIM)

for p in transformer.parameters():
    if p.dim() > 1:
        nn.init.xavier_uniform_(p)

transformer = transformer.to(DEVICE)

loss_fn = torch.nn.CrossEntropyLoss(ignore_index=PAD_IDX)

optimizer = torch.optim.Adam(transformer.parameters(), lr=0.0001, betas=(0.9, 0.98), eps=1e-9)
```

CODE IMPLEMENTATION

CONVERTING RAW SENTENCES INTO PROPER TOKENS:

```
from torch.nn.utils.rnn import pad_sequence

# helper function to club together sequential operations
def sequential_transforms(*transforms):
    def func(txt_input):
        for transform in transforms:
            txt_input = transform(txt_input)
        return txt_input
    return func

# function to add BOS/EOS and create tensor for input sequence indices
def tensor_transform(token_ids: List[int]):
    return torch.cat((torch.tensor([BOS_IDX]),
                     torch.tensor(token_ids),
                     torch.tensor([EOS_IDX])))

# src and tgt language text transforms to convert raw strings into tensors indices
text_transform = {}
for ln in [SRC_LANGUAGE, TGT_LANGUAGE]:
    text_transform[ln] = sequential_transforms(token_transform[ln], #Tokenization
                                              vocab_transform[ln], #Numericalization
                                              tensor_transform) # Add BOS/EOS and create tensor

# function to collate data samples into batch tesors
def collate_fn(batch):
    src_batch, tgt_batch = [], []
    for src_sample, tgt_sample in batch:
        src_batch.append(text_transform[SRC_LANGUAGE](src_sample.rstrip("\n")))
        tgt_batch.append(text_transform[TGT_LANGUAGE](tgt_sample.rstrip("\n")))

    src_batch = pad_sequence(src_batch, padding_value=PAD_IDX)
    tgt_batch = pad_sequence(tgt_batch, padding_value=PAD_IDX)
    return src_batch, tgt_batch
```

CODE IMPLEMENTATION

TRAINING:

```
from torch.utils.data import DataLoader

def train_epoch(model, optimizer):
    model.train()
    losses = 0
    train_iter = Multi30k(split='train', language_pair=(SRC_LANGUAGE, TGT_LANGUAGE))
    train_dataloader = DataLoader(train_iter, batch_size=BATCH_SIZE, collate_fn=collate_fn)

    for src, tgt in train_dataloader:
        src = src.to(DEVICE)
        tgt = tgt.to(DEVICE)

        tgt_input = tgt[:-1, :]

        src_mask, tgt_mask, src_padding_mask, tgt_padding_mask = create_mask(src, tgt_input)

        logits = model(src, tgt_input, src_mask, tgt_mask, src_padding_mask, tgt_padding_mask, src_padding_mask)

        optimizer.zero_grad()

        tgt_out = tgt[1:, :]
        loss = loss_fn(logits.reshape(-1, logits.shape[-1]), tgt_out.reshape(-1))
        loss.backward()

        optimizer.step()
        losses += loss.item()

    return losses / len(train_dataloader)
```

CODE IMPLEMENTATION

VALIDATION:

```
def evaluate(model):
    model.eval()
    losses = 0

    val_iter = Multi30k(split='valid', language_pair=(SRC_LANGUAGE, TGT_LANGUAGE))
    val_dataloader = DataLoader(val_iter, batch_size=BATCH_SIZE, collate_fn=collate_fn)

    for src, tgt in val_dataloader:
        src = src.to(DEVICE)
        tgt = tgt.to(DEVICE)

        tgt_input = tgt[:-1, :]

        src_mask, tgt_mask, src_padding_mask, tgt_padding_mask = create_mask(src, tgt_input)

        logits = model(src, tgt_input, src_mask, tgt_mask, src_padding_mask, tgt_padding_mask, src_padding_mask)

        tgt_out = tgt[1:, :]
        loss = loss_fn(logits.reshape(-1, logits.shape[-1]), tgt_out.reshape(-1))
        losses += loss.item()

    return losses / len(val_dataloader)
```

CODE IMPLEMENTATION

RUNNING THE MODEL:

```
from timeit import default_timer as timer
NUM_EPOCHS = 18

for epoch in range(1, NUM_EPOCHS+1):
    start_time = timer()
    train_loss = train_epoch(transformer, optimizer)
    end_time = timer()
    val_loss = evaluate(transformer)
    print(f'\tTrain Loss: {train_loss:.3f} | Train PPL: {math.exp(train_loss):7.3f}')
    print(f'\tVal. Loss: {val_loss:.3f} | Val. PPL: {math.exp(val_loss):7.3f}'')
```

CODE IMPLEMENTATION

INFERENCE :

```
def greedy_decode(model, src, src_mask, max_len, start_symbol):
    src = src.to(DEVICE)
    src_mask = src_mask.to(DEVICE)

    memory = model.encode(src, src_mask)
    ys = torch.ones(1, 1).fill_(start_symbol).type(torch.long).to(DEVICE)
    for i in range(max_len-1):
        memory = memory.to(DEVICE)
        tgt_mask = (generate_square_subsequent_mask(ys.size(0))
                    .type(torch.bool)).to(DEVICE)
        out = model.decode(ys, memory, tgt_mask)
        out = out.transpose(0, 1)
        prob = model.generator(out[:, -1])
        _, next_word = torch.max(prob, dim=1)
        next_word = next_word.item()

        ys = torch.cat([ys,
                       torch.ones(1, 1).type_as(src.data).fill_(next_word)], dim=0)
        if next_word == EOS_IDX:
            break
    return ys

# actual function to translate input sentence into target language
def translate(model: torch.nn.Module, src_sentence: str):
    model.eval()
    src = text_transform[SRC_LANGUAGE](src_sentence).view(-1, 1)
    num_tokens = src.shape[0]
    src_mask = (torch.zeros(num_tokens, num_tokens)).type(torch.bool)
    tgt_tokens = greedy_decode(
        model, src, src_mask, max_len=num_tokens + 5, start_symbol=BOS_IDX).flatten()
    return " ".join(vocab_transform[TGT_LANGUAGE].lookup_tokens(list(tgt_tokens.cpu().numpy()))).replace("<bos>", "").replace("<eos>", "")
```

CODE IMPLEMENTATION

OUR ACCURACY:
~~(100%)~~

```
print(translate(transformer, "Sie spielen Fußball."))
```

They are playing soccer .

LIMITATIONS

	Parameter Breakdown by Component:
generator	: 5,559,381 parameters
src_tok_emb	: 9,837,568 parameters
tgt_tok_emb	: 5,548,544 parameters
transformer.decoder.layers	: 7,888,896 parameters
transformer.decoder.norm	: 1,024 parameters
transformer.encoder.layers	: 4,733,952 parameters
transformer.encoder.norm	: 1,024 parameters
	Total Trainable Parameters : 33,570,389

we have roughly 30 million parameters and a data set of 50 thousand sentences.

Which is nearly 1 word for 600 parameters.

The data set happens to be too small to train a transformer with 30 million parameters which is the reason for the inaccuracy of the model.

THANK YOU