

## Import Libraries

In [1]:

```
import tensorflow as tf
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from tensorflow.keras.layers import TextVectorization
import re, string
from tensorflow.keras.layers import LSTM, Dense, Embedding, Dropout, LayerNormalization
```

In [2]:

```
df=pd.read_csv('/kaggle/input/simple-dialogs-for-chatbot/dialogs.txt', sep='\t', names=['question', 'answer'])
print(f'Dataframe size: {len(df)}')
df.head()
```

Dataframe size: 3725

Out[2]:

|   | question                            | answer                              |
|---|-------------------------------------|-------------------------------------|
| 0 | hi, how are you doing?              | i'm fine. how about yourself?       |
| 1 | i'm fine. how about yourself?       | i'm pretty good. thanks for asking. |
| 2 | i'm pretty good. thanks for asking. | no problem. so how have you been?   |

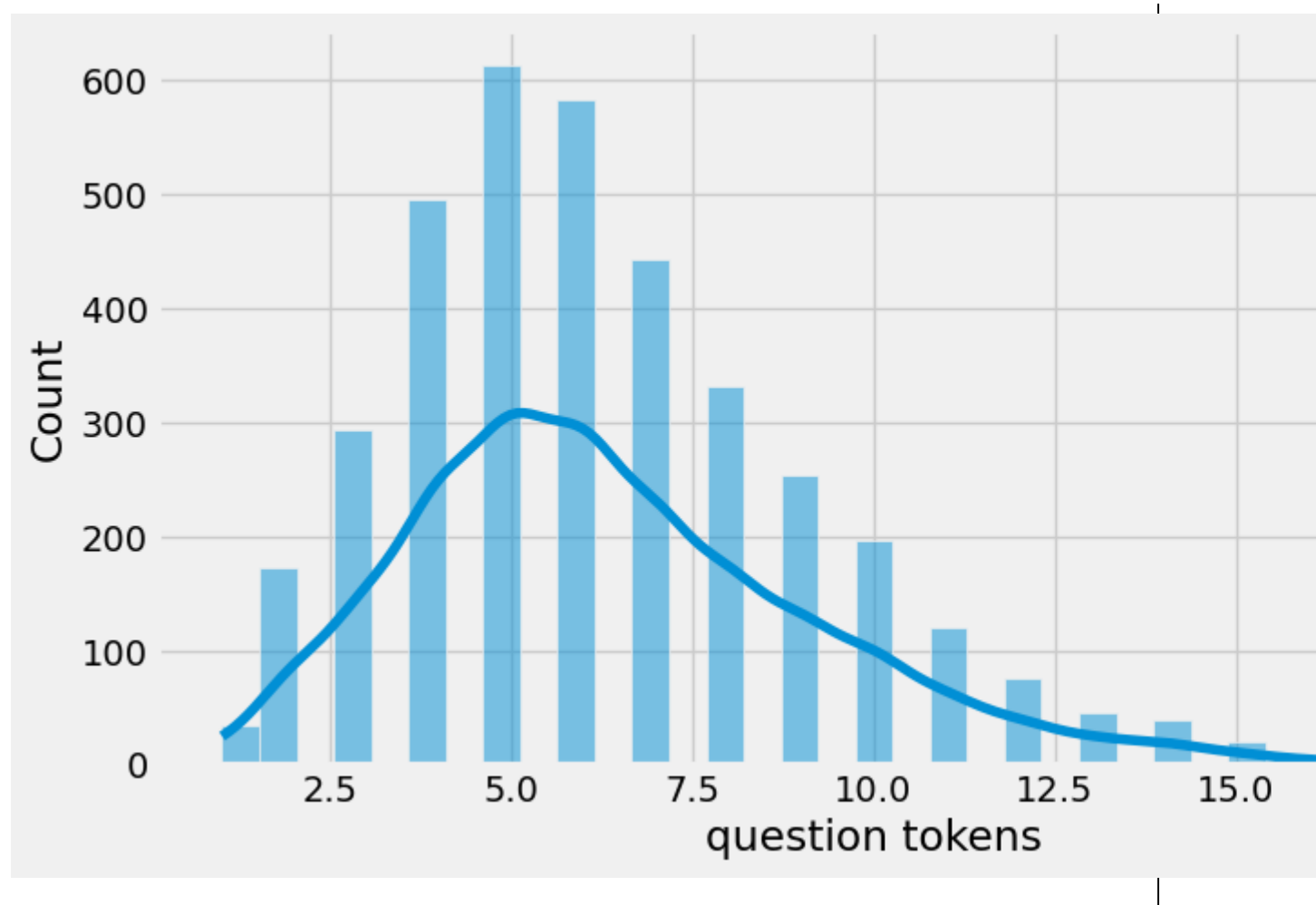
|   | question                          | answer                                   |
|---|-----------------------------------|--|
| 3 | no problem. so how have you been? | i've been great. what about you?         |
| 4 | i've been great. what about you?  | i've been good. i'm in school right now. |

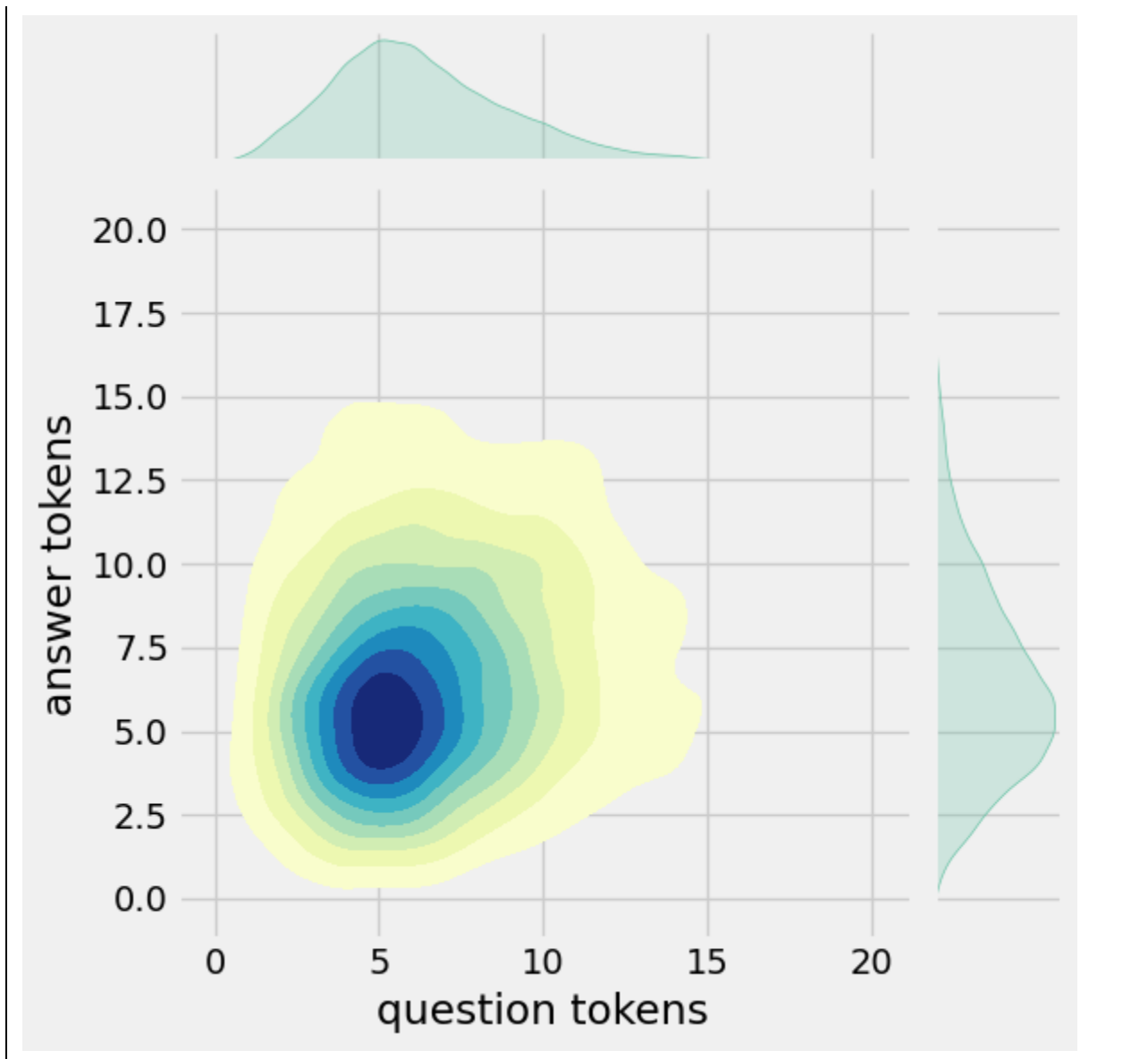
## Data Preprocessing

## Data Visualization

In [3]:

```
df['question tokens']=df['question'].apply(lambda x:len(x.split()))
df['answer tokens']=df['answer'].apply(lambda x:len(x.split()))
plt.style.use('fivethirtyeight')
fig,ax=plt.subplots(nrows=1,ncols=2,figsize=(20,5))
sns.set_palette('Set2')
sns.histplot(x=df['question tokens'],data=df,kde=True,ax=ax[0])
sns.histplot(x=df['answer tokens'],data=df,kde=True,ax=ax[1])
sns.jointplot(x='question tokens',y='answer tokens',data=df,kind='kde',fill=
True,cmap='YlGnBu')
plt.show()
```





## Text Cleaning

In [4]:

```
def clean_text(text):  
    text=re.sub('-', ' ', text.lower())  
    text=re.sub('[.]', ' . ', text)  
    text=re.sub('[1]', ' 1 ', text)  
    text=re.sub('[2]', ' 2 ', text)  
    text=re.sub('[3]', ' 3 ', text)  
    text=re.sub('[4]', ' 4 ', text)  
    text=re.sub('[5]', ' 5 ', text)  
    text=re.sub('[6]', ' 6 ', text)
```

```

text=re.sub('[7]', ' 7 ',text)
text=re.sub('[8]', ' 8 ',text)
text=re.sub('[9]', ' 9 ',text)
text=re.sub('[0]', ' 0 ',text)
text=re.sub('[,]', ' , ',text)
text=re.sub('[?]', ' ? ',text)
text=re.sub('[!]', ' ! ',text)
text=re.sub('[\$]', ' $ ',text)
text=re.sub('[&]', ' & ',text)
text=re.sub('[/]', ' / ',text)
text=re.sub('[:]', ' : ',text)
text=re.sub('[;]', ' ; ',text)
text=re.sub('[*]', ' * ',text)
text=re.sub('[\\']', ' \\' ',text)
text=re.sub('[\\"]', ' \\' ',text)
text=re.sub('\\t', ' ',text)
return text

```

```

df.drop(columns=['answer tokens','question tokens'],axis=1,inplace=True)
df['encoder_inputs']=df['question'].apply(clean_text)
df['decoder_targets']=df['answer'].apply(clean_text)+' <end>'
df['decoder_inputs']='<start> '+df['answer'].apply(clean_text)+' <end>'

```

```
df.head(10)
```

Out[4]:

|   | question                            | answer                              | encoder_inputs                          | decoder_targets                               | decoder_inputs                                   |
|---|-------------------------------------|-------------------------------------|---|---|--|
| 0 | hi, how are you doing?              | i'm fine. how about yourself?       | hi , how are you doing ?                | i ' m fine . how about yourself ? <end>       | <start> i ' m fine . how about yourself ? <end>  |
| 1 | i'm fine. how about yourself?       | i'm pretty good. thanks for asking. | i ' m fine . how about yourself ?       | i ' m pretty good . thanks for asking . <end> | <start> i ' m pretty good . thanks for asking... |
| 2 | i'm pretty good. thanks for asking. | no problem. so how have you been?   | i ' m pretty good . thanks for asking . | no problem . so how have you been ? <end>     | <start> no problem . so how have you been ? ...  |

|   | question                                 | answer                                   | encoder_inputs                                 | decoder_targets                                  | decoder_inputs                                   |
|---|--|--|--|--|--|
| 3 | no problem. so how have you been?        | i've been great. what about you?         | no problem . so how have you been ?            | i ' ve been great . what about you ? <end>       | <start> i ' ve been great . what about you ? ... |
| 4 | i've been great. what about you?         | i've been good. i'm in school right now. | i ' ve been great . what about you ?           | i ' ve been good . i ' m in school right now ... | <start> i ' ve been good . i ' m in school ri... |
| 5 | i've been good. i'm in school right now. | what school do you go to?                | i ' ve been good . i ' m in school right now . | what school do you go to ? <end>                 | <start> what school do you go to ? <end>         |
| 6 | what school do you go to?                | i go to pcc.                             | what school do you go to ?                     | i go to pcc . <end>                              | <start> i go to pcc . <end>                      |
| 7 | i go to pcc.                             | do you like it there?                    | i go to pcc .                                  | do you like it there ? <end>                     | <start> do you like it there ? <end>             |
| 8 | do you like it there?                    | it's okay. it's a really big campus.     | do you like it there ?                         | it ' s okay . it ' s a really big campus . <...> | <start> it ' s okay . it ' s a really big cam... |
| 9 | it's okay. it's a really big campus.     | good luck with school.                   | it ' s okay . it ' s a really big campus .     | good luck with school . <end>                    | <start> good luck with school . <end>            |

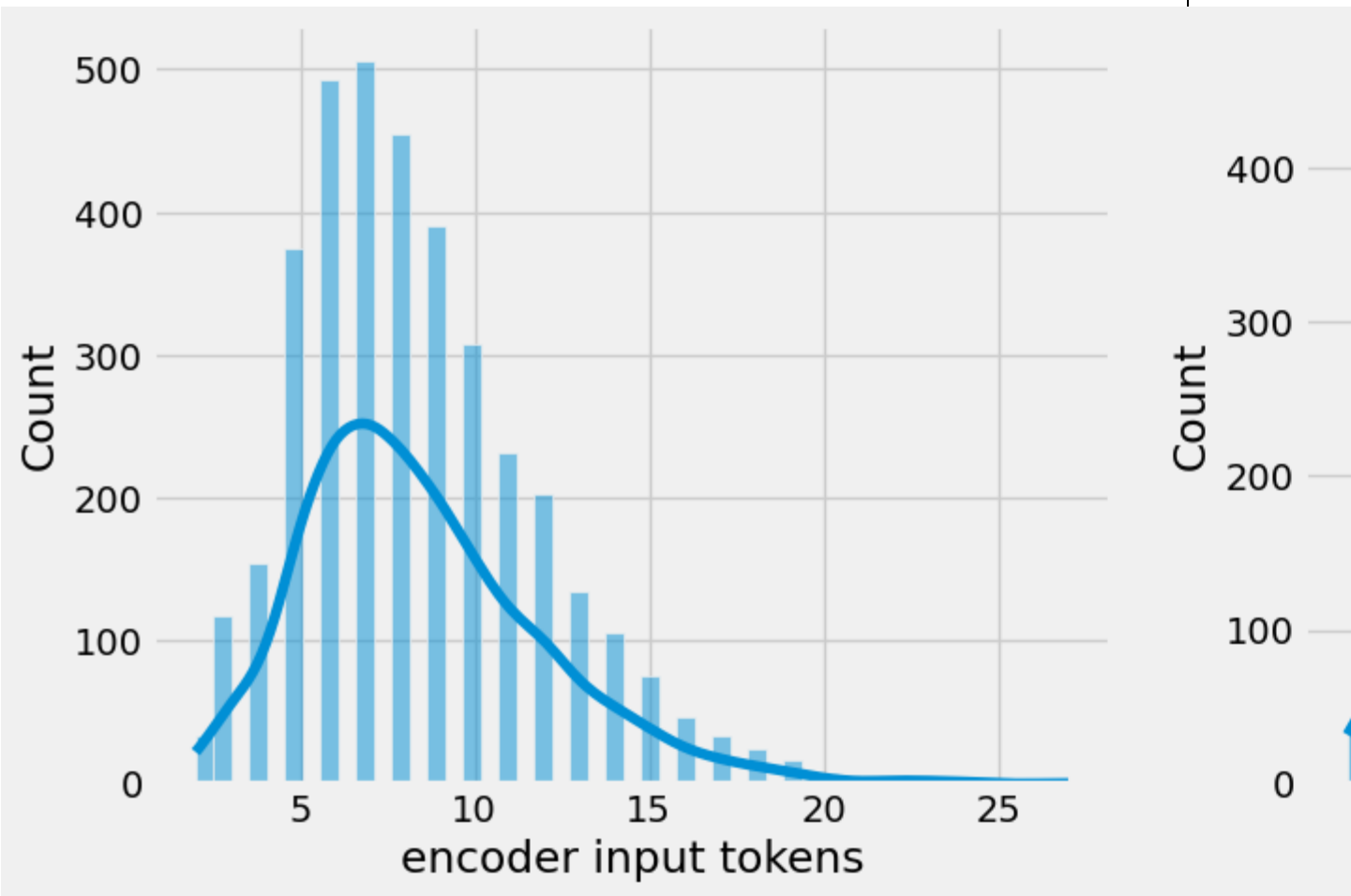
In [5]:

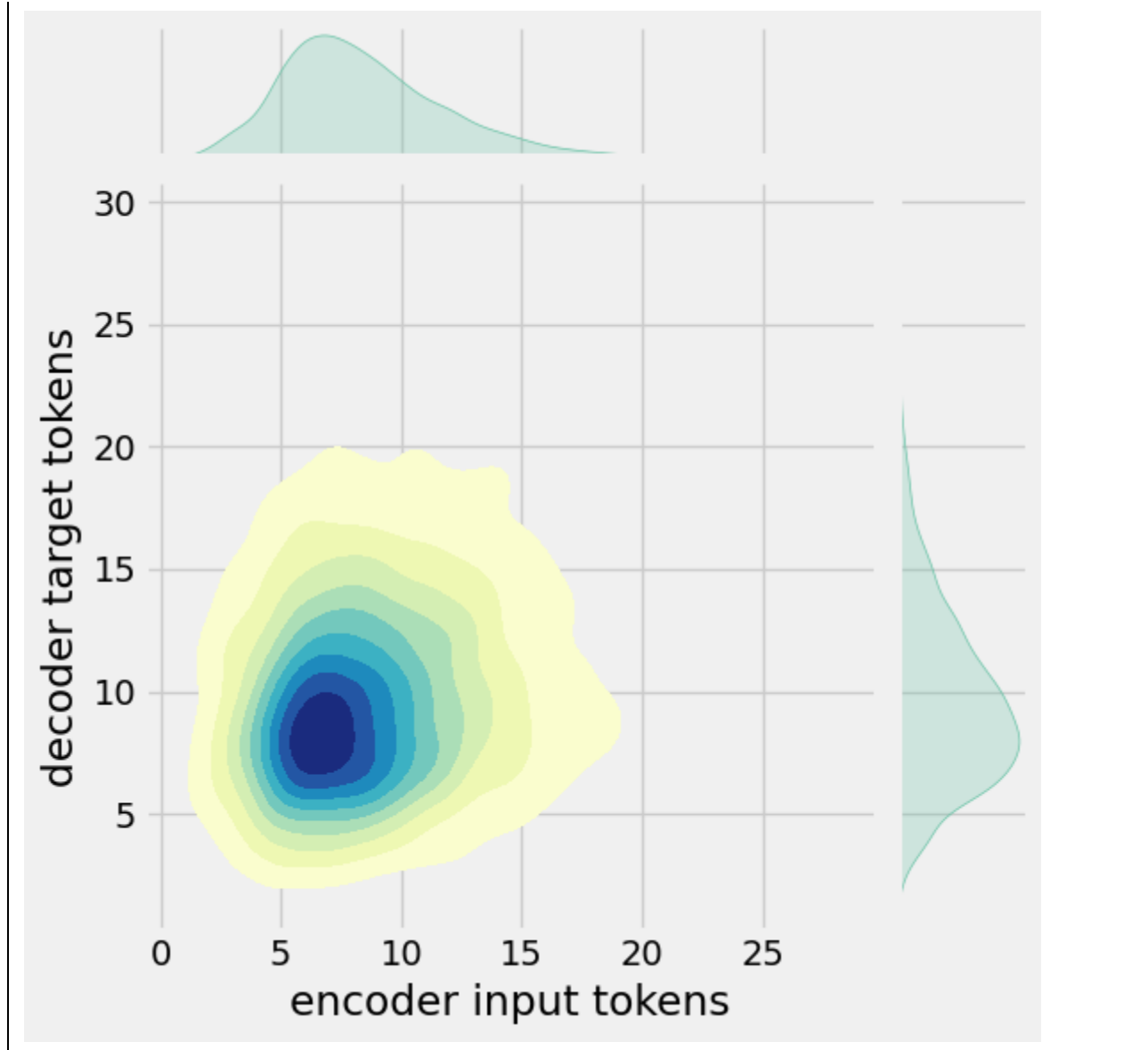
```
df['encoder input tokens']=df['encoder_inputs'].apply(lambda x:len(x.split()))
df['decoder input tokens']=df['decoder_inputs'].apply(lambda x:len(x.split()))
```

```

df['decoder target tokens']=df['decoder_targets'].apply(lambda x:len(x.split(
)))
plt.style.use('fivethirtyeight')
fig,ax=plt.subplots(nrows=1,ncols=3,figsize=(20,5))
sns.set_palette('Set2')
sns.histplot(x=df['encoder input tokens'],data=df,kde=True,ax=ax[0])
sns.histplot(x=df['decoder input tokens'],data=df,kde=True,ax=ax[1])
sns.histplot(x=df['decoder target tokens'],data=df,kde=True,ax=ax[2])
sns.jointplot(x='encoder input tokens',y='decoder target tokens',data=df,kin
d='kde',fill=True,cmap='YlGnBu')
plt.show()

```





In [6]:

```
print(f"After preprocessing: {' '.join(df[df['encoder input tokens'].max()==  
df['encoder input tokens']][['encoder_inputs'].values.tolist()})")  
print(f"Max encoder input length: {df['encoder input tokens'].max()}")  
print(f"Max decoder input length: {df['decoder input tokens'].max()}")  
print(f"Max decoder target length: {df['decoder target tokens'].max()}")  
  
df.drop(columns=['question', 'answer', 'encoder input tokens', 'decoder input t  
okens', 'decoder target tokens'], axis=1, inplace=True)  
params={  
    "vocab_size":2500,  
    "max_sequence_length":30,  
    "learning_rate":0.008,
```



```

    "batch_size":149,
    "lstm_cells":256,
    "embedding_dim":256,
    "buffer_size":10000
}
learning_rate=params['learning_rate']
batch_size=params['batch_size']
embedding_dim=params['embedding_dim']
lstm_cells=params['lstm_cells']
vocab_size=params['vocab_size']
buffer_size=params['buffer_size']
max_sequence_length=params['max_sequence_length']
df.head(10)

```

After preprocessing: for example , if your birth date is january 1 2 , 1 9 8 7 , write 0 1 / 1 2 / 8 7 .

Max encoder input length: 27

Max decoder input length: 29

Max decoder target length: 28

Out[6]:

|   | encoder_inputs                          | decoder_targets                                  | decoder_inputs                                     |
|---|---|--|--|
| 0 | hi , how are you doing ?                | i ' m fine . how about yourself ?<br><end>       | <start> i ' m fine . how about yourself ?<br><end> |
| 1 | i ' m fine . how about yourself ?       | i ' m pretty good . thanks for asking .<br><end> | <start> i ' m pretty good . thanks for asking...   |
| 2 | i ' m pretty good . thanks for asking . | no problem . so how have you been ?<br><end>     | <start> no problem . so how have you been ? ...    |
| 3 | no problem . so how have you been ?     | i ' ve been great . what about you ?<br><end>    | <start> i ' ve been great . what about you ? ...   |
| 4 | i ' ve been great . what about you ?    | i ' ve been good . i ' m in school right now ... | <start> i ' ve been good . i ' m in school ri...   |

|   | encoder_inputs                                 | decoder_targets                                  | decoder_inputs                                   |
|---|--|--|--|
| 5 | i ' ve been good . i ' m in school right now . | what school do you go to ? <end>                 | <start> what school do you go to ? <end>         |
| 6 | what school do you go to ?                     | i go to pcc . <end>                              | <start> i go to pcc . <end>                      |
| 7 | i go to pcc .                                  | do you like it there ? <end>                     | <start> do you like it there ? <end>             |
| 8 | do you like it there ?                         | it ' s okay . it ' s a really big campus . <...> | <start> it ' s okay . it ' s a really big cam... |
| 9 | it ' s okay . it ' s a really big campus .     | good luck with school . <end>                    | <start> good luck with school . <end>            |

## Tokenization

In [7]:

```
vectorize_layer=TextVectorization(
    max_tokens=vocab_size,
    standardize=None,
    output_mode='int',
    output_sequence_length=max_sequence_length
)
vectorize_layer.adapt(df['encoder_inputs']+' '+df['decoder_targets']+' <start> <end>')
vocab_size=len(vectorize_layer.get_vocabulary())
print(f'Vocab size: {len(vectorize_layer.get_vocabulary())}')
print(f'{vectorize_layer.get_vocabulary()[:12]}')
```

Vocab size: 2443

['', '[UNK]', '<end>', '.', '<start>', '"', 'i', '?', 'you', ',', 'the', 'to']

In [8]:

```

def sequences2ids(sequence):
    return vectorize_layer(sequence)

def ids2sequences(ids):
    decode=''
    if type(ids)==int:
        ids=[ids]
    for id in ids:
        decode+=vectorize_layer.get_vocabulary()[id]+' '
    return decode

x=sequences2ids(df['encoder_inputs'])
yd=sequences2ids(df['decoder_inputs'])
y=sequences2ids(df['decoder_targets'])

print(f'Question sentence: hi , how are you ?')
print(f'Question to tokens: {sequences2ids("hi , how are you ?")[:10]}')
print(f'Encoder input shape: {x.shape}')
print(f'Decoder input shape: {yd.shape}')
print(f'Decoder target shape: {y.shape}')

```

```

Question sentence: hi , how are you ?
Question to tokens: [1971    9   45   24    8    7    0    0    0    0]
Encoder input shape: (3725, 30)
Decoder input shape: (3725, 30)
Decoder target shape: (3725, 30)

```

In [9]:

```

print(f'Encoder input: {x[0][:12]} ...')
print(f'Decoder input: {yd[0][:12]} ...')    # shifted by one time step of the
print(f'Decoder target: {y[0][:12]} ...')    target as input to decoder is the output of the previous timestep

```

```

Encoder input: [1971    9   45   24    8  194    7    0    0    0    0
0] ...
Decoder input: [  4   6   5  38 646   3  45  41 563   7   2   0] ...
Decoder target: [  6   5  38 646   3  45  41 563   7   2   0   0] ...

```

In [10]:

```

data=tf.data.Dataset.from_tensor_slices((x,yd,y))
data=data.shuffle(buffer_size)

train_data=data.take(int(.9*len(data)))
train_data=train_data.cache()
train_data=train_data.shuffle(buffer_size)
train_data=train_data.batch(batch_size)

```

```

train_data=train_data.prefetch(tf.data.AUTOTUNE)
train_data_iterator=train_data.as_numpy_iterator()

val_data=data.skip(int(.9*len(data))).take(int(.1*len(data)))
val_data=val_data.batch(batch_size)
val_data=val_data.prefetch(tf.data.AUTOTUNE)

_=train_data_iterator.next()
print(f'Number of train batches: {len(train_data)}')
print(f'Number of training data: {len(train_data)*batch_size}')
print(f'Number of validation batches: {len(val_data)}')
print(f'Number of validation data: {len(val_data)*batch_size}')
print(f'Encoder Input shape (with batches): {_[0].shape}')
print(f'Decoder Input shape (with batches): {_[1].shape}')
print(f'Target Output shape (with batches): {_[2].shape}')

```

```

Number of train batches: 23
Number of training data: 3427
Number of validation batches: 3
Number of validation data: 447
Encoder Input shape (with batches): (149, 30)
Decoder Input shape (with batches): (149, 30)
Target Output shape (with batches): (149, 30)

```

## Build Models

### Build Encoder

In [ ]:

linkcode

In [11]:

```

class Encoder(tf.keras.models.Model):
    def __init__(self, units, embedding_dim, vocab_size, *args, **kwargs) -> None
    :
        super().__init__(*args, **kwargs)
        self.units=units
        self.vocab_size=vocab_size
        self.embedding_dim=embedding_dim
        self.embedding=Embedding(
            vocab_size,
            embedding_dim,

```

```

        name='encoder_embedding',
        mask_zero=True,
        embeddings_initializer=tf.keras.initializers.GlorotNormal()
    )
    self.normalize=LayerNormalization()
    self.lstm=LSTM(
        units,
        dropout=.4,
        return_state=True,
        return_sequences=True,
        name='encoder_lstm',
        kernel_initializer=tf.keras.initializers.GlorotNormal()
    )

def call(self,encoder_inputs):
    self.inputs=encoder_inputs
    x=self.embedding(encoder_inputs)
    x=self.normalize(x)
    x=Dropout(.4)(x)
    encoder_outputs,encoder_state_h,encoder_state_c=self.lstm(x)
    self.outputs=[encoder_state_h,encoder_state_c]
    return encoder_state_h,encoder_state_c

encoder=Encoder(lstm_cells,embedding_dim,vocab_size,name='encoder')
encoder.call(_[0])

```

Out[11]:

```

(<tf.Tensor: shape=(149, 256), dtype=float32, numpy=
array([[ 0.16966951, -0.10419625, -0.12700348, ..., -0.12251794,
        0.10568858,  0.14841646],
       [ 0.08443093,  0.08849293, -0.09065959, ..., -0.00959182,
        0.10152507, -0.12077457],
       [ 0.03628462, -0.02653611, -0.11506603, ..., -0.14669597,
        0.10292757,  0.13625325],
       ...,
       [-0.14210635, -0.12942064, -0.03288083, ...,  0.0568463 ,
        -0.02598592, -0.22455114],
       [ 0.20819993,  0.01196991, -0.09635217, ..., -0.18782297,
        0.10233591,  0.20114912],
       [ 0.1164271 , -0.07769038, -0.06414707, ..., -0.06539135,
        -0.05518465,  0.25142196]], dtype=float32)>,
<tf.Tensor: shape=(149, 256), dtype=float32, numpy=
array([[ 0.34589   , -0.30134732, -0.43572   , ..., -0.3102559 ,
        0.34630865,  0.2613009 ],
       [ 0.14154069,  0.17045322, -0.17749965, ..., -0.02712595,
        0.17292541, -0.2922624 ],
       [ 0.07106856, -0.0739173 , -0.3641197 , ..., -0.3794833 ,

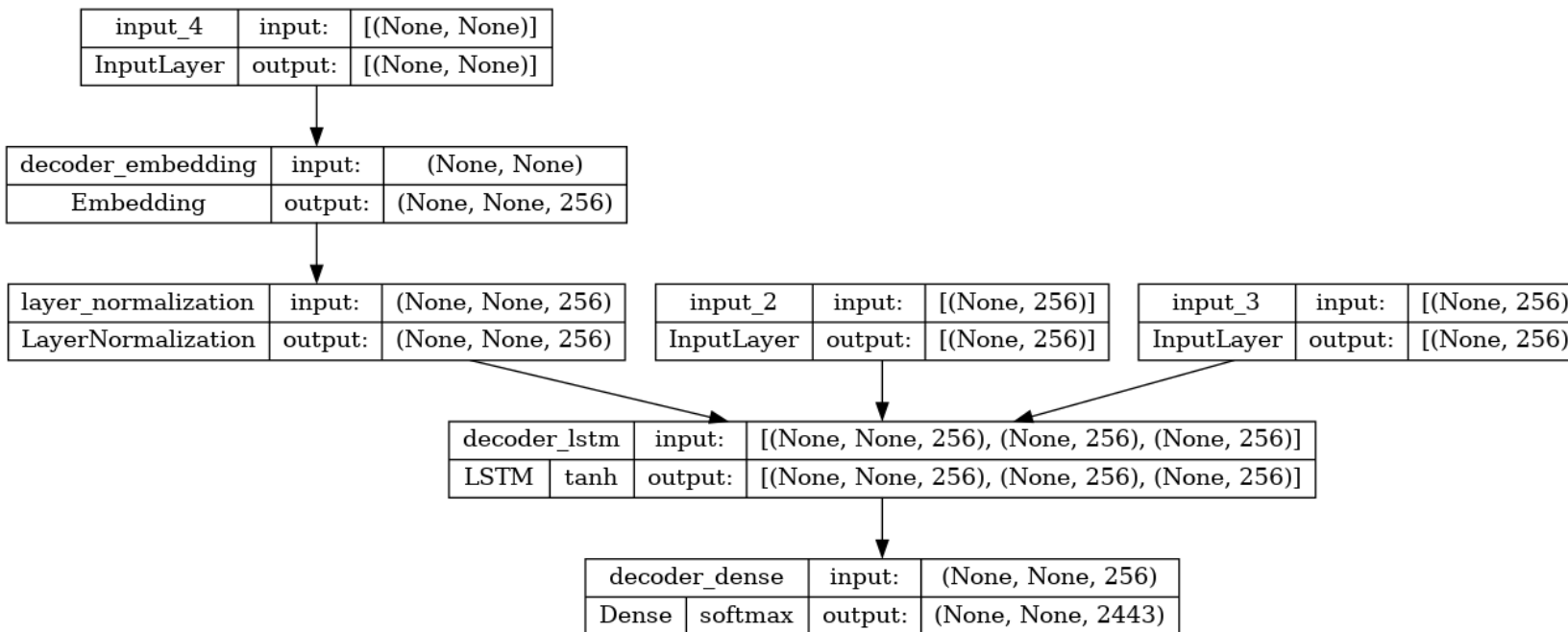
```

```

0.36470377, 0.23766585],
...,
[-0.2582597, -0.25323495, -0.06649272, ..., 0.16527973,
-0.04292646, -0.58768904],
[ 0.43155715, 0.03135502, -0.33463806, ..., -0.47625306,
0.33486888, 0.35035062],
[ 0.23173636, -0.20141824, -0.22034441, ..., -0.16035017,
-0.17478186, 0.48899865]], dtype=float32)>)

```

## Create Inference Model



## Time to Chat

```
Def print_conversation(texts):
```

```
    For text in texts:
```

```
        Print(f'You: {text}')
```

```
        Print(f'Bot: {chatbot(text)}')
```

```
        Print('=====')
```

```
Print_conversation([
```

```
    'hi',
```

```
    'do yo know me?',
```

```
    'what is your name?',
```

```
    'you are bot?',
```

'hi, how are you doing?',  
"I'm pretty good. Thanks for asking.",  
"Don't ever be in a hurry",  
""I'm gonna put some dirt in your eye """,  
""You're trash """,  
""I've read all your research on nano-technology """,  
""You want forgiveness? Get religion""",  
""While you're using the bathroom, I'll order some food.""",  
""Wow! That's terrible.""",  
""We'll be here forever.""",  
""I need something that's reliable.""",  
""A speeding car ran a red light, killing the girl.""",  
""Tomorrow we'll have rice and fish for lunch.""",  
""I like this restaurant because they give you free bread.""

])

You: hi

Bot: I have to go to the bathroom.

---

You: do yo know me?

Bot: yes, it's too close to the other.

---

You: what is your name?

Bot: I have to walk the house.

---

You: you are bot?

Bot: no, I have. All my life.

---

You: hi, how are you doing?

Bot: I'm going to be a teacher.

---

You: I'm pretty good. Thanks for asking.

Bot: no problem. I'll have to give you the english assignments from my mind.

---

You: Don't ever be in a hurry

Bot: it's not a great.

---

You: I'm gonna put some dirt in your eye

Bot: that's a good idea.

---

You: You're trash

Bot: the tv news is reporting a bank robbery.

---

You: I've read all your research on nano-technology

Bot: it's the weather. I've gone around the world.

---

You: You want forgiveness? Get religion

Bot: no, I'll be my.

---

You: While you're using the bathroom, I'll order some food.

Bot: don't order for me. I've been a cheater.

---

You: Wow! That's terrible.

Bot: never park your car under the house.

---

You: We'll be here forever.

Bot: we'll be there in half an hour.

---

You: I need something that's reliable.

Bot: you need a car with low mileage.

---

You: A speeding car ran a red light, killing the girl.

Bot: what happened?

---

You: Tomorrow we'll have rice and fish for lunch.

Bot: I'll make a sandwich.

---

You: I like this restaurant because they give you free bread.

Bot: well, I think that's a good idea.

---