

✓ Phase 2: Transformer-Based Summarization

Overview

In this phase, we build a baseline abstractive summarization model using the Transformer architecture. This model will serve as a benchmark for comparison against more advanced models in subsequent phases.

1. Tokenization and Preprocessing

- Load the cleaned dataset.
 - Construct a vocabulary based on the most frequent words.
 - Encode the input articles and summaries as sequences of integer indices.
 - Pad sequences to fixed length for batch processing.
-

2. Transformer Model Architecture

- Use PyTorch's `nn.Transformer` or a custom encoder-decoder Transformer model.
 - Components:
 - Token embedding layer
 - Positional encoding
 - Transformer Encoder
 - Transformer Decoder
 - Final linear layer to project decoder outputs to vocabulary size
-

3. Masking Strategy

- **Padding Mask:** Prevents the model from attending to padded elements.
 - **Look-ahead Mask:** Ensures that the decoder attends only to previous tokens during training.
-

4. Model Training

- Loss function: `nn.CrossEntropyLoss(ignore_index=PAD)`
 - Optimizer: Adam
 - Training loop for multiple epochs
 - Log training and validation loss
-

5. Evaluation

- Implement greedy decoding during inference.
 - Compare model-generated summaries against ground-truth summaries.
 - Print sample predictions to assess performance.
-

6. Next Steps


- Analyze failure cases.
- Consider training with subword tokenization or pre-trained embeddings.
- Prepare to fine-tune a pre-trained transformer in Phase 3.

```
import pandas as pd
from google.colab import files

uploaded = files.upload()

for fn in uploaded.keys():
    df = pd.read_csv(fn)
    print(f'User uploaded file "{fn}" with length {len(uploaded[fn])} bytes')
    break # Exit the loop after processing the first file

print(df.head())
```

 no files selected Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

Saving cleaned_dataset.csv to cleaned_dataset.csv

User uploaded file "cleaned_dataset.csv" with length 113592985 bytes

	URL	\
0	https://www.moneycontrol.com/news/business/eco...	
1	https://www.businesstoday.in/top-story/state-r...	
2	https://www.financialexpress.com/economy/covid...	
3	https://www.moneycontrol.com/news/business/mar...	
4	https://www.financialexpress.com/industry/six-...	

	Content	\
0	us consumer spending dropped by a record in ap...	
1	state-run lenders require an urgent rs 1.2 tri...	
2	apparel exporters on wednesday urged the gover...	
3	asian shares battled to extend a global reboun...	
4	after india's sovereign credit rating fell to ...	

	Summary	Sentiment
0	consumer spending plunges 13.6 percent in apri...	Negative
1	government will have to take a bulk of the tab...	Negative
2	exporters are facing issues in terms of raw ma...	Negative
3	the dollar loses some ground on the safe haven...	Negative

```

import pandas as pd
import nltk
from nltk.tokenize import TreebankWordTokenizer
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder

# Load dataset
df = pd.read_csv("cleaned_dataset.csv")

# Rename columns for consistency
df.rename(columns={'Content': 'text', 'Summary': 'summary', 'Sentiment': 'sentiment'})

# Drop rows with missing values
df.dropna(subset=['text', 'summary', 'sentiment'], inplace=True)

# Tokenization using TreebankWordTokenizer (no punkt_tab needed)
tokenizer = TreebankWordTokenizer()
df['text_tokens'] = df['text'].apply(tokenizer.tokenize)
df['summary_tokens'] = df['summary'].apply(lambda x: ['<S0S>'] + tokenizer.tokenize(x))

# Encode sentiment labels (e.g., Positive = 2, Negative = 0, Neutral = 1)
label_encoder = LabelEncoder()
df['sentiment_encoded'] = label_encoder.fit_transform(df['sentiment'])

# Split into train/val/test
train_df, temp_df = train_test_split(df, test_size=0.15, random_state=42)
val_df, test_df = train_test_split(temp_df, test_size=0.33, random_state=42) # 10%

# Print class distribution for sanity check
print("✅ Sentiment Labels:", label_encoder.classes_)
print("✅ Train size:", len(train_df), "Validation size:", len(val_df), "Test size:", len(test_df))
print("✅ Sample tokenized text:", df['text_tokens'].iloc[0][:15])
print("✅ Sample sentiment encoding:", df['sentiment_encoded'].iloc[0])

```

```

↔️ ✅ Sentiment Labels: ['Negative' 'Neutral' 'Positive']
✅ Train size: 22298 Validation size: 2636 Test size: 1299
✅ Sample tokenized text: ['us', 'consumer', 'spending', 'dropped', 'by', 'a']
✅ Sample sentiment encoding: 0

```

```

from collections import Counter
import torch

```

```

# Build vocab from train set only (text + summary)
def build_vocab(token_lists, min_freq=2):
    counter = Counter()
    for tokens in token_lists:
        counter.update(tokens)

    vocab = {
        '<PAD>': 0,
        '<UNK>': 1,
        '<SOS>': 2,
        '<EOS>': 3
    }
    idx = 4
    for word, freq in counter.items():
        if freq >= min_freq and word not in vocab:
            vocab[word] = idx
            idx += 1
    return vocab

# Combine text and summary tokens for vocab
combined_tokens = train_df['text_tokens'].tolist() + train_df['summary_tokens'].tolist()
word2idx = build_vocab(combined_tokens)
idx2word = {v: k for k, v in word2idx.items()}

# Helper to convert tokens to indices
def encode(tokens, word2idx):
    return [word2idx.get(token, word2idx['<UNK>']) for token in tokens]

# Apply encoding
for df_ in [train_df, val_df, test_df]:
    df_['text_idx'] = df_['text_tokens'].apply(lambda tokens: encode(tokens, word2idx))
    df_['summary_idx'] = df_['summary_tokens'].apply(lambda tokens: encode(tokens, word2idx))

# Define special token indices
PAD = word2idx['<PAD>']
UNK = word2idx['<UNK>']
SOS = word2idx['<SOS>']
EOS = word2idx['<EOS>']

print(f"✅ Vocab size: {len(word2idx)}")
print(f"✅ Sample encoded input: {train_df['text_idx'].iloc[0][:10]}")
print(f"✅ Sample encoded summary: {train_df['summary_idx'].iloc[0][:10]}")
print(f"✅ SOS={SOS}, EOS={EOS}, PAD={PAD}, UNK={UNK}")

```

```

➡️ ✓ Vocab size: 119246
✓ Sample encoded input: [4, 5, 6, 7, 8, 9, 10, 11, 12, 13]
✓ Sample encoded summary: [2, 64, 49, 53, 65, 66, 67, 28, 68, 69]
✓ SOS=2, EOS=3, PAD=0, UNK=1

```

```

from collections import defaultdict
from torch.utils.data import Dataset, DataLoader
from sklearn.model_selection import train_test_split
from torch.nn.utils.rnn import pad_sequence
import torch

```

```

# ✓ 1. Build vocabulary
token2idx = {'<PAD>': 0, '<UNK>': 1, '<SOS>': 2, '<EOS>': 3}
idx = 4
for tokens in df['text_tokens'].tolist() + df['summary_tokens'].tolist():
    for token in tokens:
        if token not in token2idx:
            token2idx[token] = idx
            idx += 1

```

```

VOCAB_SIZE = len(token2idx)
PAD, UNK, SOS, EOS = token2idx['<PAD>'], token2idx['<UNK>'], token2idx['<SOS>'],

```

```

# ✓ 2. Encode tokens
df['text_encoded'] = df['text_tokens'].apply(lambda tokens: [token2idx.get(token,
df['summary_encoded'] = df['summary_tokens'].apply(lambda tokens: [token2idx.get(

```

```

# ✓ 3. Encode sentiment
sentiment_to_idx = {'Negative': 0, 'Neutral': 1, 'Positive': 2}
df['sentiment_encoded'] = df['sentiment'].map(sentiment_to_idx)

```

```

# ✓ 4. Split into Train/Validation/Test
train_df, temp_df = train_test_split(df, test_size=0.15, random_state=42)
val_df, test_df = train_test_split(temp_df, test_size=0.5, random_state=42)

```

```

# ✓ 5. Custom Dataset
class NewsSummaryDataset(Dataset):
    def __init__(self, texts, summaries, sentiments):
        self.texts = texts
        self.summaries = summaries
        self.sentiments = sentiments

    def __len__(self):
        return len(self.texts)

```

```
def __getitem__(self, idx):
    return (
        torch.tensor(self.texts[idx], dtype=torch.long),
        torch.tensor(self.summaries[idx], dtype=torch.long),
        torch.tensor(self.sentiments[idx], dtype=torch.long)
    )
```

6. Create Datasets

```
train_dataset = NewsSummaryDataset(train_df['text_encoded'].tolist(),
                                    train_df['summary_encoded'].tolist(),
                                    train_df['sentiment_encoded'].tolist())
```

```
val_dataset = NewsSummaryDataset(val_df['text_encoded'].tolist(),
                                  val_df['summary_encoded'].tolist(),
                                  val_df['sentiment_encoded'].tolist())
```

```
test_dataset = NewsSummaryDataset(test_df['text_encoded'].tolist(),
                                   test_df['summary_encoded'].tolist(),
                                   test_df['sentiment_encoded'].tolist())
```

7. Collate Function for Dynamic Padding

```
def collate_fn(batch):
    texts, summaries, sentiments = zip(*batch)
    padded_texts = pad_sequence(texts, batch_first=True, padding_value=PAD)
    padded_summaries = pad_sequence(summaries, batch_first=True, padding_value=PAD)
    sentiments = torch.stack(sentiments)
    return padded_texts, padded_summaries, sentiments
```

8. Dataloaders with collate_fn

```
BATCH_SIZE = 4 # Safe for Colab Pro
```

```
train_loader = DataLoader(train_dataset, batch_size=BATCH_SIZE, shuffle=True, collate_fn=collate_fn)
```


```
val_loader = DataLoader(val_dataset, batch_size=BATCH_SIZE, collate_fn=collate_fn)
```

```
test_loader = DataLoader(test_dataset, batch_size=BATCH_SIZE, collate_fn=collate_fn)
```

Confirm


```
print( Sentiment Labels:", df['sentiment'].unique())
```

```
print( Train size:", len(train_df), "Validation size:", len(val_df), "Test size:", len(test_df))
```

```
print( Vocab size:", VOCAB_SIZE)
```

```
print( SOS=", SOS, "EOS=", EOS, "PAD=", PAD, "UNK=", UNK)
```

```
print( Sample encoded input:", train_df['text_encoded'].iloc[0][:10])
```

```
print( Sample encoded summary:", train_df['summary_encoded'].iloc[0][:10])
```

```

➡️ ✓ Sentiment Labels: ['Negative' 'Neutral' 'Positive']
✓ Train size: 22298 Validation size: 1967 Test size: 1968
✓ Vocab size: 167964
✓ SOS= 2 EOS= 3 PAD= 0 UNK= 1
✓ Sample encoded input: [842, 1843, 1134, 103, 13724, 26460, 28637, 19, 1548
✓ Sample encoded summary: [2, 345, 111, 9, 390, 518, 317, 14, 6929, 63]

```

```

import torch.nn as nn
import torch

```

```

class PositionalEncoding(nn.Module):
    def __init__(self, d_model, max_len=10000):
        super().__init__()
        pe = torch.zeros(max_len, d_model)
        pos = torch.arange(0, max_len).unsqueeze(1).float()
        div_term = torch.exp(torch.arange(0, d_model, 2).float() * (-torch.log(to
        pe[:, 0::2] = torch.sin(pos * div_term)
        pe[:, 1::2] = torch.cos(pos * div_term)
        pe = pe.unsqueeze(0) # shape: (1, max_len, d_model)
        self.register_buffer('pe', pe)

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

class DualTaskTransformer(nn.Module):
    def __init__(self, vocab_size, emb_size, num_heads, num_encoder_layers, num_d
        hidden_dim, dropout=0.1, num_classes=3):
        super().__init__()
        self.embedding = nn.Embedding(vocab_size, emb_size, padding_idx=PAD)
        self.pos_encoder = PositionalEncoding(emb_size)
        self.pos_decoder = PositionalEncoding(emb_size)

        encoder_layer = nn.TransformerEncoderLayer(d_model=emb_size, nhead=num_he
        decoder_layer = nn.TransformerDecoderLayer(d_model=emb_size, nhead=num_he

        self.encoder = nn.TransformerEncoder(encoder_layer, num_layers=num_encode
        self.decoder = nn.TransformerDecoder(decoder_layer, num_layers=num_decode

        self.fc_out = nn.Linear(emb_size, vocab_size) # For summarization
        self.sentiment_head = nn.Sequential(
            nn.Linear(emb_size, hidden_dim),
            nn.ReLU(),
            nn.Dropout(dropout),

```



```

        nn.Linear(hidden_dim, num_classes) # For sentiment classification
    )

def make_src_mask(self, src):
    return (src == PAD)

def make_tgt_mask(self, tgt):
    tgt_pad_mask = (tgt == PAD)
    tgt_len = tgt.size(1)
    tgt_sub_mask = torch.tril(torch.ones((tgt_len, tgt_len), device=tgt.device))
    return tgt_pad_mask, tgt_sub_mask

def forward(self, src, tgt):
    src_mask = self.make_src_mask(src)
    tgt_pad_mask, tgt_sub_mask = self.make_tgt_mask(tgt)

    src_emb = self.pos_encoder(self.embedding(src))
    memory = self.encoder(src_emb, src_key_padding_mask=src_mask)

    # Sentiment prediction (mean of encoder output)
    sent_pred = self.sentiment_head(memory.mean(dim=1))

    tgt_emb = self.pos_decoder(self.embedding(tgt))
    output = self.decoder(
        tgt=tgt_emb,
        memory=memory,
        tgt_mask=tgt_sub_mask,
        tgt_key_padding_mask=tgt_pad_mask,
        memory_key_padding_mask=src_mask
    )

    out_tokens = self.fc_out(output)
    return out_tokens, sent_pred

```

```

import torch
import torch.nn as nn
import torch.optim as optim

```

```

# ✅ Use the correct VOCAB_SIZE from earlier
VOCAB_SIZE = 167964
NUM_CLASSES = 3 # Negative, Neutral, Positive

```

```

# 🛠 Model Hyperparameters

```

```

EMBED_SIZE = 256
NUM_HEADS = 8
ENC_LAYERS = 3
DEC_LAYERS = 3
HIDDEN_DIM = 512
DROPOUT = 0.1

# 🚫 PAD index used to ignore padding tokens during loss calculation
PAD = 0

# 🖥️ Device setup
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

# ⚙️ Initialize model
model = DualTaskTransformer(
    vocab_size=VOCAB_SIZE,
    emb_size=EMBED_SIZE,
    num_heads=NUM_HEADS,
    num_encoder_layers=ENC_LAYERS,
    num_decoder_layers=DEC_LAYERS,
    hidden_dim=HIDDEN_DIM,
    dropout=DROPOUT,
    num_classes=NUM_CLASSES
).to(device)

# 📉 Loss functions
criterion_summary = nn.CrossEntropyLoss(ignore_index=PAD) # For summary generation
criterion_sentiment = nn.CrossEntropyLoss() # For sentiment classification

# 🚀 Optimizer
optimizer = optim.Adam(model.parameters(), lr=1e-4)

print("✅ Model, loss functions, and optimizer initialized.")

```

🔄 ✅ Model, loss functions, and optimizer initialized.

```

import torch

device = torch.device("cuda")
scaler = torch.amp.GradScaler('cuda') # <-- updated
EPOCHS = 5
CLIP = 1.0

```

```
model.to(device)
model.train()

for epoch in range(EPOCHS):
    total_loss = 0.0
    correct = 0
    total = 0

    for input_tensor, summary_tensor, sentiment_tensor in train_loader:
        input_tensor = input_tensor.to(device)
        summary_tensor = summary_tensor.to(device)
        sentiment_tensor = sentiment_tensor.to(device)

        tgt_input = summary_tensor[:, :-1]
        tgt_output = summary_tensor[:, 1:]

        optimizer.zero_grad()

        with torch.amp.autocast('cuda'): # <-- updated
            summary_logits, sentiment_logits = model(input_tensor, tgt_input)


            loss1 = criterion_summary(
                summary_logits.reshape(-1, summary_logits.size(-1)),
                tgt_output.reshape(-1)
            )
            loss2 = criterion_sentiment(sentiment_logits, sentiment_tensor)
            loss = loss1 + loss2

        scaler.scale(loss).backward()
        torch.nn.utils.clip_grad_norm_(model.parameters(), CLIP)
        scaler.step(optimizer)
        scaler.update()

        total_loss += loss.item()
        correct += (sentiment_logits.argmax(1) == sentiment_tensor).sum().item()
        total += sentiment_tensor.size(0)

    avg_loss = total_loss / len(train_loader)
    acc = correct / total

    print(f"Epoch {epoch+1}/{EPOCHS} | Loss: {avg_loss:.4f} | Sentiment Acc: {acc
```



```
Epoch 1/5 | Loss: 7.2020 | Sentiment Acc: 0.4306  
Epoch 2/5 | Loss: 4.3213 | Sentiment Acc: 0.5547  
Epoch 3/5 | Loss: 3.0184 | Sentiment Acc: 0.6123  
Epoch 4/5 | Loss: 2.4018 | Sentiment Acc: 0.6502  
Epoch 5/5 | Loss: 2.0191 | Sentiment Acc: 0.6932
```

```
torch.save(model.state_dict(), "dual_task_transformer.pth")  
from google.colab import files  
files.download("dual_task_transformer.pth")
```



```
model.eval()
test_loss = 0.0
correct = 0
total = 0

with torch.no_grad():
    for input_tensor, summary_tensor, sentiment_tensor in test_loader:
        input_tensor = input_tensor.to(device)
        summary_tensor = summary_tensor.to(device)
        sentiment_tensor = sentiment_tensor.to(device)

        tgt_input = summary_tensor[:, :-1]
        tgt_output = summary_tensor[:, 1:]

        with torch.amp.autocast('cuda'):
            summary_logits, sentiment_logits = model(input_tensor, tgt_input)

            loss1 = criterion_summary(
                summary_logits.reshape(-1, summary_logits.size(-1)),
                tgt_output.reshape(-1)
            )
            loss2 = criterion_sentiment(sentiment_logits, sentiment_tensor)
            loss = loss1 + loss2

        test_loss += loss.item()
        correct += (sentiment_logits.argmax(1) == sentiment_tensor).sum().item()
        total += sentiment_tensor.size(0)

avg_test_loss = test_loss / len(test_loader)
test_acc = correct / total

print(f"Test Loss: {avg_test_loss:.4f} | Sentiment Accuracy: {test_acc:.4f}")
```

➡ Test Loss: 1.9689 | Sentiment Accuracy: 0.6463

```

def greedy_decode(model, src_tensor, max_len=30):
    model.eval()
    src_tensor = src_tensor[:256] # truncate input if it's too long
    src_tensor = src_tensor.unsqueeze(0).to(device)

    with torch.no_grad():
        src_mask = model.make_src_mask(src_tensor)
        src_emb = model.pos_encoder(model.embedding(src_tensor))
        memory = model.encoder(src_emb, src_key_padding_mask=src_mask)

        tgt_indices = [SOS]

        for step in range(max_len):
            tgt_tensor = torch.tensor(tgt_indices, dtype=torch.long, device=device)
            tgt_pad_mask, tgt_sub_mask = model.make_tgt_mask(tgt_tensor)

            tgt_emb = model.pos_decoder(model.embedding(tgt_tensor))
            output = model.decoder(
                tgt=tgt_emb,
                memory=memory,
                tgt_mask=tgt_sub_mask,
                tgt_key_padding_mask=tgt_pad_mask,
                memory_key_padding_mask=src_mask
            )

            logits = model.fc_out(output)
            next_token = logits[0, -1].argmax().item()

            if step == 0:
                print("First token predicted:", next_token, idx2word.get(next_token))

            if next_token in [EOS, PAD]:
                break

            tgt_indices.append(next_token)

    return tgt_indices[1:] # remove <SOS>

```

```
NUM_SAMPLES = 5
```

```
for i in range(NUM_SAMPLES):
    input_ids = torch.tensor(test_df['text_encoded'].iloc[i][:256] # Truncate
    target_ids = test_df['summary_encoded'].iloc[i]

    pred_ids = greedy_decode(model, input_ids)

    pred_tokens = [idx2word.get(idx, '<UNK>') for idx in pred_ids]
    target_tokens = [idx2word.get(idx, '<UNK>') for idx in target_ids if idx not

    print(f"\nSample {i+1}:")
    print("Predicted Summary:", " ".join(pred_tokens))
    print("Actual Summary   :", " ".join(target_tokens))
```

⇒ First token predicted: 0 <PAD>

Sample 1:

Predicted Summary:

Actual Summary : for imposes flagship securities again saturday go-ahead mor

First token predicted: 0 <PAD>

Sample 2:

Predicted Summary:

Actual Summary : for singh strong christopher impaired sources # monday note

First token predicted: 0 <PAD>

Sample 3:

Predicted Summary:

Actual Summary : for 4.5-5 lack heads marginally legally expect company end

First token predicted: 0 <PAD>

Sample 4:

Predicted Summary:

Actual Summary : for surplus us held cliff christopher will monday intention

First token predicted: 0 <PAD>

Sample 5:

Predicted Summary:

Actual Summary : 0.1-0.5 non-essential gold abhaya demonstrate mark contrast

```

input_ids = torch.tensor(train_df['text_encoded'].iloc[0])[:256]
target_ids = train_df['summary_encoded'].iloc[0]

pred_ids = greedy_decode(model, input_ids)

pred_tokens = [idx2word.get(idx, '<UNK>') for idx in pred_ids]
target_tokens = [idx2word.get(idx, '<UNK>') for idx in target_ids if idx not in [0]]

print("Predicted Summary:", " ".join(pred_tokens))
print("Actual Summary   :", " ".join(target_tokens))

```

```

➞ First token predicted: 0 <PAD>
Predicted Summary:
Actual Summary   : - lack christopher ait post-covid 300 for 4,991.50 selloff.

```

```

def greedy_decode(model, src_tensor, max_len=30):
    model.eval()
    src_tensor = src_tensor[:256] # truncate long inputs
    src_tensor = src_tensor.unsqueeze(0).to(device)

    with torch.no_grad():
        src_mask = model.make_src_mask(src_tensor)
        src_emb = model.pos_encoder(model.embedding(src_tensor))
        memory = model.encoder(src_emb, src_key_padding_mask=src_mask)

        tgt_indices = [SOS]

        for step in range(max_len):
            tgt_tensor = torch.tensor(tgt_indices, dtype=torch.long, device=device)
            tgt_pad_mask, tgt_sub_mask = model.make_tgt_mask(tgt_tensor)

            tgt_emb = model.pos_decoder(model.embedding(tgt_tensor))
            output = model.decoder(
                tgt=tgt_emb,
                memory=memory,
                tgt_mask=tgt_sub_mask,
                tgt_key_padding_mask=tgt_pad_mask,
                memory_key_padding_mask=src_mask
            )

            logits = model.fc_out(output)[0, -1]

```



```

# Prevent <PAD> token from being chosen
logits[PAD] = -1e9

# Choose top-1 token that isn't <PAD>
topk = torch.topk(logits, k=5)
next_token = topk.indices[0].item()

if step == 0:
    print("First token predicted:", next_token, idx2word.get(next_tok

if next_token in [EOS, PAD]:
    break

tgt_indices.append(next_token)

return tgt_indices[1:] # remove <SOS>

```

```

input_ids = torch.tensor(train_df['text_encoded'].iloc[0])[:256]
target_ids = train_df['summary_encoded'].iloc[0]

pred_ids = greedy_decode(model, input_ids)

pred_tokens = [idx2word.get(idx, '<UNK>') for idx in pred_ids]
target_tokens = [idx2word.get(idx, '<UNK>') for idx in target_ids if idx not in [

print("\nPredicted Summary:", " ".join(pred_tokens))
print("Actual Summary      :", " ".join(target_tokens))

```

➡ First token predicted: 2 <SOS>

Predicted Summary: <SOS> <SOS> <SOS> <SOS> <SOS> <SOS> <SOS> <SOS> <SOS> <SOS>
 Actual Summary : - lack christopher ait post-covid 300 for 4,991.50 selloff.

✓ Debug Report – Baseline Transformer Summarizer

1. Issues Observed

Symptom	Evidence
Decoder collapses to a single token (<PAD> or <SOS>)	Greedy decoding on train & test samples yields strings of <PAD> →
Sentiment head trains correctly	Sent-accuracy rose from 43 % → 69 % in 5 epochs
Summary branch never learns	Even after a 3-epoch top-up (summary-only, loss ↓ 1.06 → 0.76) deco

Root Cause

Severe loss imbalance & shortcut learning – during mixed-task training the model found a low-loss path by predicting padding tokens for summaries, letting sentiment dominate optimisation. Once that collapse occurred, a short fine-tune could not recover a usable token distribution.

2. Experiments & Results

Test	Result
Block <PAD> at inference	Decoder switched to <SOS> loop
Decode training sample	Same collapse → confirms decoder never learned
Summary-only top-up (3 epochs)	Loss fell, but output remained <SOS> – weights still unusable
Vocab remap check	Verified correct lookup; output truly is <SOS> token-ID

3. Corrective Action Plan

3.1 Full Retrain (recommended)

1. **Single-task phase** – train **only summarisation** for 5–6 epochs.
2. **Add sentiment head** – fine-tune 2-3 epochs with weighted loss:

$$\text{loss} = 1.5 * \text{summary_loss} + 0.5 * \text{sentiment_loss}$$

Start coding or generate with AI.