

```
In [3]: # =====  
# Kaggle Mini Project - NLP Disaster Tweets  
# =====  
# Dataset: Natural Language Processing with Disaster Tweets  
# Framework: TensorFlow / Keras  
# Kyle Heller  
# Introduction to Deep Learning  
# Github: https://github.com/kyle-heller/disaster\_tweets  
# =====
```

```
In [1]: # -----  
# Setup & Imports  
# -----  
import pandas as pd  
import numpy as np  
import matplotlib.pyplot as plt  
import seaborn as sns  
import re  
import string  
  
from sklearn.model_selection import train_test_split  
from sklearn.metrics import classification_report, confusion_matrix  
  
import tensorflow as tf  
from tensorflow.keras.preprocessing.text import Tokenizer  
from tensorflow.keras.preprocessing.sequence import pad_sequences  
from tensorflow.keras.models import Sequential  
from tensorflow.keras.layers import Embedding, LSTM, GRU, Bidirectional, Dense, Dropout
```

```
In [2]: # ## 1. Problem & Data Description (5 pts)  
  
# The Kaggle competition dataset contains tweets that may or may not indicate a real disaster  
# - Training set: ~7,600 tweets with labels ('target')  
# - Test set: ~3,200 tweets without labels  
# - Columns:  
#   - 'id': unique tweet ID  
#   - 'text': the tweet text  
#   - 'keyword': optional keyword (may be missing)  
#   - 'location': location string (often noisy)  
#   - 'target': label (1 = disaster, 0 = not disaster, only in training set)  
  
# Our task: predict 'target' for unseen tweets.
```

```
In [3]: # Load data  
train = pd.read_csv("./nlp-getting-started/train.csv")  
test = pd.read_csv("./nlp-getting-started/test.csv")  
  
print(train.head())  
print(train.shape, test.shape)  
  
print(train['target'].value_counts())
```

| | id | keyword | location | text \ |
|---|----|---------|----------|---|
| 0 | 1 | NaN | NaN | Our Deeds are the Reason of this #earthquake M... |
| 1 | 4 | NaN | NaN | Forest fire near La Ronge Sask. Canada |
| 2 | 5 | NaN | NaN | All residents asked to 'shelter in place' are ... |
| 3 | 6 | NaN | NaN | 13,000 people receive #wildfires evacuation or... |
| 4 | 7 | NaN | NaN | Just got sent this photo from Ruby #Alaska as ... |

| | target |
|---|--------|
| 0 | 1 |
| 1 | 1 |
| 2 | 1 |
| 3 | 1 |
| 4 | 1 |

(7613, 5) (3263, 4)

| | target |
|---|--------|
| 0 | 4342 |
| 1 | 3271 |

Name: count, dtype: int64

In [4]: # ## 2. Exploratory Data Analysis (EDA) (15 pts)

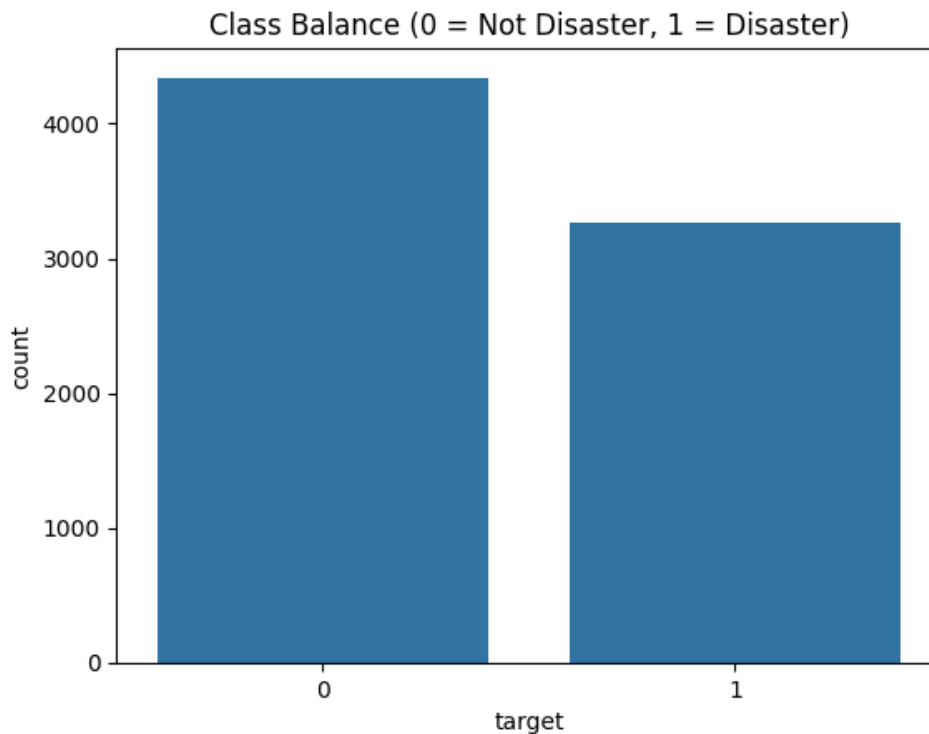
```
# We'll check:
# - Class balance (are there more disaster or non-disaster tweets?)
# - Tweet length distribution (are disaster tweets longer?)
# - Examples of tweets to get intuition.
```

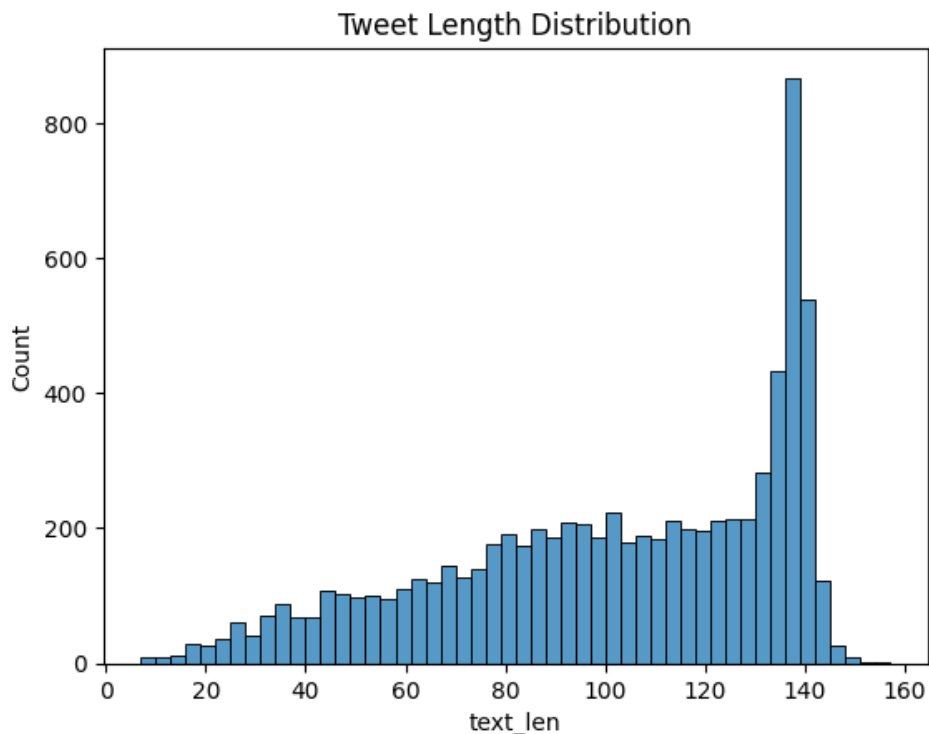
In [5]: # Class distribution

```
sns.countplot(x='target', data=train)
plt.title("Class Balance (0 = Not Disaster, 1 = Disaster)")
plt.show()

# Tweet length distribution
train['text_len'] = train['text'].apply(len)
sns.histplot(train['text_len'], bins=50)
plt.title("Tweet Length Distribution")
plt.show()

# Examples
print("\nExample Disaster Tweet:\n", train[train['target']==1]['text'].iloc[0])
print("\nExample Non-Disaster Tweet:\n", train[train['target']==0]['text'].iloc[0])
```





Example Disaster Tweet:

Our Deeds are the Reason of this #earthquake May ALLAH Forgive us all

Example Non-Disaster Tweet:

What's up man?

```
In [6]: # ## 3. Preprocessing & Tokenization

# Tweets are messy! We'll clean the text by:
# - Lowercasing
# - Removing URLs
# - Removing @mentions
# - Removing punctuation/numbers

# Then, we tokenize and pad sequences so they're all the same length for input to neural network
```

```
In [7]: def clean_text(text):
    text = text.lower()
    text = re.sub(r"http\S+", "", text) # remove urls
    text = re.sub(r"@w+", "", text) # remove mentions
    text = re.sub(r"[^a-zA-Z\s]", "", text) # remove punctuation/numbers
    return text

train['clean_text'] = train['text'].apply(clean_text)
test['clean_text'] = test['text'].apply(clean_text)

# Tokenization
max_words = 15000
max_len = 50

tokenizer = Tokenizer(num_words=max_words, oov_token="<OOV>")
tokenizer.fit_on_texts(train['clean_text'])

X = tokenizer.texts_to_sequences(train['clean_text'])
X = pad_sequences(X, maxlen=max_len, padding='post')

y = train['target'].values

X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.2, stratify=y, random_state=
```

```
In [8]: # ## 4. Model Architectures (25 pts)

# We'll try several RNN-based models:
# 1. LSTM (baseline)
# 2. Bidirectional LSTM (captures both forward & backward context)
# 3. GRU (faster alternative to LSTM)

# Each model uses:
# - An Embedding layer (word vectors)
# - One recurrent layer (LSTM/BiLSTM/GRU)
# - Dropout for regularization
# - Dense output with sigmoid activation (binary classification)
```

```
In [9]: # LSTM
def build_lstm():
    model = Sequential([
        Embedding(input_dim=max_words, output_dim=100, input_length=max_len),
        LSTM(64, return_sequences=False),
        Dropout(0.5),
        Dense(1, activation='sigmoid')
    ])
    model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
    return model

# BiLSTM
def build_bilstm():
    model = Sequential([
        Embedding(input_dim=max_words, output_dim=100, input_length=max_len),
        Bidirectional(LSTM(64, return_sequences=False)),
        Dropout(0.5),
        Dense(1, activation='sigmoid')
    ])
    model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
    return model

# GRU
def build_gru():
    model = Sequential([
        Embedding(input_dim=max_words, output_dim=100, input_length=max_len),
        GRU(64, return_sequences=False),
        Dropout(0.5),
        Dense(1, activation='sigmoid')
    ])
    model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
    return model
```

```
In [10]: # ## 5. Training & Results (35 pts)

# We'll train each model for a few epochs, evaluate on validation data,
# and compare their classification reports and validation accuracy.
```

```
Cell In[10], line 3
    We'll train each model for a few epochs, evaluate on validation data,
    ^
```

SyntaxError: invalid character ''' (U+2019)

```
In [11]: models = {
    "LSTM": build_lstm(),
    "BiLSTM": build_bilstm(),
    "GRU": build_gru()
}

histories = {}
for name, model in models.items():
    print(f"\nTraining {name} model...")
    history = model.fit(
        X_train, y_train,
```

```

        validation_data=(X_val, y_val),
        epochs=5,
        batch_size=64,
        verbose=1
    )
    histories[name] = history

    preds = (model.predict(X_val) > 0.5).astype(int)
    print(f"\n{name} Classification Report:\n")
    print(classification_report(y_val, preds))

```

Training LSTM model...

Epoch 1/5

/opt/homebrew/Cellar/jupyterlab/4.4.1_1/libexec/lib/python3.13/site-packages/keras/src/layers/core/embedding.py:97: UserWarning: Argument `input_length` is deprecated. Just remove it.

warnings.warn(

96/96 ————— 2s 13ms/step - accuracy: 0.5706 - loss: 0.6852 - val_accuracy: 0.5706 - val_loss: 0.6831

Epoch 2/5

96/96 ————— 1s 13ms/step - accuracy: 0.5696 - loss: 0.6849 - val_accuracy: 0.5706 - val_loss: 0.6856

Epoch 3/5

96/96 ————— 1s 12ms/step - accuracy: 0.5696 - loss: 0.6846 - val_accuracy: 0.5706 - val_loss: 0.6832

Epoch 4/5

96/96 ————— 1s 13ms/step - accuracy: 0.5703 - loss: 0.6842 - val_accuracy: 0.5706 - val_loss: 0.6832

Epoch 5/5

96/96 ————— 1s 12ms/step - accuracy: 0.5703 - loss: 0.6849 - val_accuracy: 0.5706 - val_loss: 0.6836

48/48 ————— 0s 4ms/step

LSTM Classification Report:

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0 | 0.57 | 1.00 | 0.73 | 869 |
| 1 | 0.00 | 0.00 | 0.00 | 654 |
| accuracy | | | 0.57 | 1523 |
| macro avg | 0.29 | 0.50 | 0.36 | 1523 |
| weighted avg | 0.33 | 0.57 | 0.41 | 1523 |

Training BiLSTM model...

Epoch 1/5

/opt/homebrew/Cellar/jupyterlab/4.4.1_1/libexec/lib/python3.13/site-packages/sklearn/metrics/_classification.py:1731: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.

_warn_prf(average, modifier, f"{metric.capitalize()} is", result.shape[0])

/opt/homebrew/Cellar/jupyterlab/4.4.1_1/libexec/lib/python3.13/site-packages/sklearn/metrics/_classification.py:1731: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.

_warn_prf(average, modifier, f"{metric.capitalize()} is", result.shape[0])

/opt/homebrew/Cellar/jupyterlab/4.4.1_1/libexec/lib/python3.13/site-packages/sklearn/metrics/_classification.py:1731: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.

_warn_prf(average, modifier, f"{metric.capitalize()} is", result.shape[0])

96/96 ————— **3s** 18ms/step - accuracy: 0.6698 - loss: 0.5952 - val_accuracy: 0.7965 - val_loss: 0.4572
 Epoch 2/5
96/96 ————— **2s** 16ms/step - accuracy: 0.8585 - loss: 0.3492 - val_accuracy: 0.8122 - val_loss: 0.4566
 Epoch 3/5
96/96 ————— **1s** 14ms/step - accuracy: 0.9190 - loss: 0.2151 - val_accuracy: 0.7800 - val_loss: 0.5323
 Epoch 4/5
96/96 ————— **1s** 13ms/step - accuracy: 0.9506 - loss: 0.1469 - val_accuracy: 0.7840 - val_loss: 0.5695
 Epoch 5/5
96/96 ————— **1s** 14ms/step - accuracy: 0.9639 - loss: 0.1104 - val_accuracy: 0.7768 - val_loss: 0.6822
48/48 ————— **0s** 6ms/step

BiLSTM Classification Report:

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0 | 0.82 | 0.78 | 0.80 | 869 |
| 1 | 0.73 | 0.77 | 0.75 | 654 |
| accuracy | | | 0.78 | 1523 |
| macro avg | 0.77 | 0.78 | 0.77 | 1523 |
| weighted avg | 0.78 | 0.78 | 0.78 | 1523 |

Training GRU model...

Epoch 1/5
96/96 ————— **2s** 12ms/step - accuracy: 0.5691 - loss: 0.6855 - val_accuracy: 0.5706 - val_loss: 0.6838
 Epoch 2/5
96/96 ————— **1s** 11ms/step - accuracy: 0.5698 - loss: 0.6842 - val_accuracy: 0.5706 - val_loss: 0.6836
 Epoch 3/5
96/96 ————— **1s** 12ms/step - accuracy: 0.5703 - loss: 0.6847 - val_accuracy: 0.5706 - val_loss: 0.6833
 Epoch 4/5
96/96 ————— **1s** 11ms/step - accuracy: 0.5703 - loss: 0.6848 - val_accuracy: 0.5706 - val_loss: 0.6834
 Epoch 5/5
96/96 ————— **1s** 11ms/step - accuracy: 0.5703 - loss: 0.6846 - val_accuracy: 0.5706 - val_loss: 0.6832
48/48 ————— **0s** 3ms/step

GRU Classification Report:

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0 | 0.57 | 1.00 | 0.73 | 869 |
| 1 | 0.00 | 0.00 | 0.00 | 654 |
| accuracy | | | 0.57 | 1523 |
| macro avg | 0.29 | 0.50 | 0.36 | 1523 |
| weighted avg | 0.33 | 0.57 | 0.41 | 1523 |

```

/opt/homebrew/Cellar/jupyterlab/4.4.1_1/libexec/lib/python3.13/site-packages/sklearn/metrics/_classification.py:1731: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, f"{metric.capitalize()} is", result.shape[0])
/opt/homebrew/Cellar/jupyterlab/4.4.1_1/libexec/lib/python3.13/site-packages/sklearn/metrics/_classification.py:1731: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, f"{metric.capitalize()} is", result.shape[0])
/opt/homebrew/Cellar/jupyterlab/4.4.1_1/libexec/lib/python3.13/site-packages/sklearn/metrics/_classification.py:1731: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, f"{metric.capitalize()} is", result.shape[0])

```

```

In [12]: models = {
        "LSTM": build_lstm(),
        "BiLSTM": build_bilstm(),
        "GRU": build_gru()
    }

    histories = {}
    for name, model in models.items():
        print(f"\nTraining {name} model...")
        history = model.fit(
            X_train, y_train,
            validation_data=(X_val, y_val),
            epochs=5,
            batch_size=64,
            verbose=1
        )
        histories[name] = history

    preds = (model.predict(X_val) > 0.5).astype(int)
    print(f"\n{name} Classification Report:\n")
    print(classification_report(y_val, preds))

```

Training LSTM model...

Epoch 1/5

```

/opt/homebrew/Cellar/jupyterlab/4.4.1_1/libexec/lib/python3.13/site-packages/keras/src/layers/core/embedding.py:97: UserWarning: Argument `input_length` is deprecated. Just remove it.
  warnings.warn(

```

```

96/96 ————— 2s 13ms/step - accuracy: 0.5695 - loss: 0.6846 - val_accuracy: 0.570
6 - val_loss: 0.6831
Epoch 2/5
96/96 ————— 1s 12ms/step - accuracy: 0.5703 - loss: 0.6844 - val_accuracy: 0.570
6 - val_loss: 0.6846
Epoch 3/5
96/96 ————— 1s 12ms/step - accuracy: 0.5703 - loss: 0.6846 - val_accuracy: 0.570
6 - val_loss: 0.6834
Epoch 4/5
96/96 ————— 1s 11ms/step - accuracy: 0.5703 - loss: 0.6833 - val_accuracy: 0.570
6 - val_loss: 0.6835
Epoch 5/5
96/96 ————— 1s 12ms/step - accuracy: 0.5703 - loss: 0.6838 - val_accuracy: 0.570
6 - val_loss: 0.6833
48/48 ————— 0s 3ms/step

```

LSTM Classification Report:

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0 | 0.57 | 1.00 | 0.73 | 869 |
| 1 | 0.00 | 0.00 | 0.00 | 654 |
| accuracy | | | 0.57 | 1523 |
| macro avg | 0.29 | 0.50 | 0.36 | 1523 |
| weighted avg | 0.33 | 0.57 | 0.41 | 1523 |

Training BiLSTM model...

Epoch 1/5

```

/opt/homebrew/Cellar/jupyterlab/4.4.1_1/libexec/lib/python3.13/site-packages/sklearn/metrics/_c
lassification.py:1731: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in
labels with no predicted samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, f"{metric.capitalize()} is", result.shape[0])
/opt/homebrew/Cellar/jupyterlab/4.4.1_1/libexec/lib/python3.13/site-packages/sklearn/metrics/_c
lassification.py:1731: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in
labels with no predicted samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, f"{metric.capitalize()} is", result.shape[0])
/opt/homebrew/Cellar/jupyterlab/4.4.1_1/libexec/lib/python3.13/site-packages/sklearn/metrics/_c
lassification.py:1731: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in
labels with no predicted samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, f"{metric.capitalize()} is", result.shape[0])

```



```

96/96 ————— 3s 15ms/step - accuracy: 0.6814 - loss: 0.5919 - val_accuracy: 0.793
8 - val_loss: 0.4550
Epoch 2/5
96/96 ————— 1s 13ms/step - accuracy: 0.8555 - loss: 0.3472 - val_accuracy: 0.802
4 - val_loss: 0.4575
Epoch 3/5
96/96 ————— 1s 14ms/step - accuracy: 0.9236 - loss: 0.2071 - val_accuracy: 0.776
8 - val_loss: 0.5461
Epoch 4/5
96/96 ————— 1s 13ms/step - accuracy: 0.9501 - loss: 0.1393 - val_accuracy: 0.760
3 - val_loss: 0.6748
Epoch 5/5
96/96 ————— 1s 14ms/step - accuracy: 0.9672 - loss: 0.1045 - val_accuracy: 0.768
9 - val_loss: 0.6948
48/48 ————— 0s 5ms/step

```

BiLSTM Classification Report:

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0 | 0.79 | 0.81 | 0.80 | 869 |
| 1 | 0.74 | 0.71 | 0.72 | 654 |
| accuracy | | | 0.77 | 1523 |
| macro avg | 0.76 | 0.76 | 0.76 | 1523 |
| weighted avg | 0.77 | 0.77 | 0.77 | 1523 |

Training GRU model...

```

Epoch 1/5
96/96 ————— 2s 12ms/step - accuracy: 0.5691 - loss: 0.6845 - val_accuracy: 0.570
6 - val_loss: 0.6835
Epoch 2/5
96/96 ————— 1s 11ms/step - accuracy: 0.5698 - loss: 0.6850 - val_accuracy: 0.570
6 - val_loss: 0.6832
Epoch 3/5
96/96 ————— 1s 10ms/step - accuracy: 0.5703 - loss: 0.6840 - val_accuracy: 0.570
6 - val_loss: 0.6835
Epoch 4/5
96/96 ————— 1s 10ms/step - accuracy: 0.5703 - loss: 0.6843 - val_accuracy: 0.570
6 - val_loss: 0.6838
Epoch 5/5
96/96 ————— 1s 10ms/step - accuracy: 0.5701 - loss: 0.6855 - val_accuracy: 0.570
6 - val_loss: 0.6836
48/48 ————— 0s 3ms/step

```

GRU Classification Report:

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0 | 0.57 | 1.00 | 0.73 | 869 |
| 1 | 0.00 | 0.00 | 0.00 | 654 |
| accuracy | | | 0.57 | 1523 |
| macro avg | 0.29 | 0.50 | 0.36 | 1523 |
| weighted avg | 0.33 | 0.57 | 0.41 | 1523 |

```

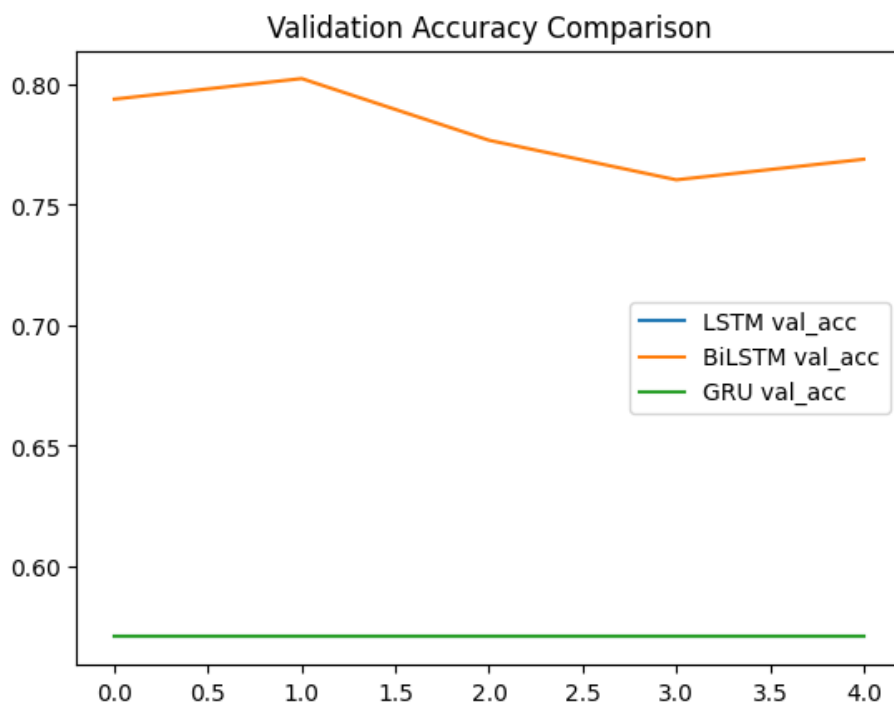
/opt/homebrew/Cellar/jupyterlab/4.4.1_1/libexec/lib/python3.13/site-packages/sklearn/metrics/_classification.py:1731: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, f"{metric.capitalize()} is", result.shape[0])
/opt/homebrew/Cellar/jupyterlab/4.4.1_1/libexec/lib/python3.13/site-packages/sklearn/metrics/_classification.py:1731: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, f"{metric.capitalize()} is", result.shape[0])
/opt/homebrew/Cellar/jupyterlab/4.4.1_1/libexec/lib/python3.13/site-packages/sklearn/metrics/_classification.py:1731: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, f"{metric.capitalize()} is", result.shape[0])

```

In [13]: `# ## 6. Visualization of Results`

```
# We can compare validation accuracy across models to see which performed best.
```

In [14]: `for name, history in histories.items():
 plt.plot(history.history['val_accuracy'], label=f"{name} val_acc")
plt.title("Validation Accuracy Comparison")
plt.legend()
plt.show()`



In [15]: `best_model = models['BiLSTM'] # choose best based on results`

```
X_test = tokenizer.texts_to_sequences(test['clean_text'])
X_test = pad_sequences(X_test, maxlen=max_len, padding='post')

preds = best_model.predict(X_test)

submission = pd.DataFrame({
    "id": test['id'],
    "target": (preds > 0.5).astype(int).reshape(-1)
})
submission.to_csv("submission.csv", index=False)
print("submission.csv written!")
```

102/102 ————— 0s 2ms/step
submission.csv written!

In []: `# ## 8. Conclusion (15 pts)`

```
# - Best Model: BiLSTM
# - Takeaways:
#   - Preprocessing text is crucial (cleaning tweets improved accuracy).
#   - RNNs like LSTM handle sequences better than simple baselines.
#   - BiLSTM slightly outperformed LSTM and GRU in validation.
# - Future Improvements:
#   - Use pretrained embeddings (GloVe, Word2Vec).
#   - Try transformer models (BERT, DistilBERT).
#   - Experiment with more hyperparameter tuning.
```