**Python for Data Science and AI using Hugging Face Pretrained Models**

Isariya Triampetch

Pichaya Detchwongchaikul

(Group 4)

Tokyo International University

Python for Data Science and AI

Mr. Rafik Hamza

December 7th, 2023

# Introduction

What is a Hugging Face? According to Maning (2021), "Hugging Face is an open-source platform that provides tools and resources for working on natural language processing (NLP) and computer vision projects." This shows that Hugging Face is an artificial intelligence research, development business, and open-source community. Its primary goal is to give users access to all the necessary tools, libraries, and resources so they can work on natural language process (NLP) models for their own advantage. An example used of pretrained models is to test summarizing systems, which can automatically create summaries of lengthy texts, and are being developed using trained models. Nallapati (2016) used a pretrained model named BART to train a text summarization system that can generate summaries of documents in different genres. In Addition, hugging face is famous for its use in Chatbots and Conversational AI. It is highlighted that Hugging Face is an invaluable open-source platform that empowers users to develop and deploy effective natural language processing (NLP) solutions.

## How can pretrained models accelerate data science and AI?

By providing several key benefits that streamline the development process and enhance the capability of AI solutions such as reduced development time, Improved accuracy, and increased accessibility. Pretrained models save data scientists and artificial intelligence developers a great deal of time and effort by removing the requirement to train models from scratch (Fu, 2022). According to this, researchers can swiftly fine-tune a pretrained model for their specific tasks, which leads them to focus on more complex and innovative aspects of their projects, such as feature engineering and model evaluation. Meanwhile, pretrained models frequently produce innovative outcomes on a range of natural language processing tasks, offering a solid foundation for customization to meet requirements. This leads to reducing time of

building high-performing models for data scientists' experimentation (Devlin, 2019). Lastly, for people with no prior experience with natural language processing (NLP), finding, downloading, and using pretrained models is simple with Hugging Face. As a result, natural language processing  (NLP) is now more approachable and democratic for academics and developers (Wolf, 2019). In summary, pre trained models accelerate the creation and use of natural language processing (NLP) based solutions, making them essential tools for data scientists and artificial intelligence engineers.

## Problem and Motivation

**The importance of pretrained model to data Science and Artificial Intelligence**

Hugging Face pretrained models are an effective means of utilizing the most recent advancements in natural language processing (NLP) technology while customizing the model for particular applications or domains. This leads to the fact that Hugging Face to be one of the significant equipment for data scientists. As stated above, pretrained models provide many benefits to users such as, reduced development time, and improved accuracy. According to Han et al. (2021) article said through the process of fine-tuning, artificial intelligence practitioners can develop highly performant models with limited task-specific data by leveraging the massive knowledge gained from large-scale data.

**Why is this problem important to data science and AI?**

1. **Difficulty of training NLP models**

The challenge of training a natural language processing (NLP) model is recognized as difficult since it demands a substantial amount of time and computational power. On the

contrary, the procedure has been considerably streamlined by the availability of pertain models. Now that these models have been trained on huge datasets, developers can use them to fine-tune them for particular purposes (Hugging Face, n.d.). With no need to start from scratch with a model, this method enables the development of customized natural language processing (NLP) applications, saving time and effort while maintaining excellent performance and accuracy levels.

## 2. Lack of labeled data

It is difficult to train natural language processing (NLP) models due to the lack of labeled data, which is needed to educate the models how to correctly comprehend and interpret human language. It can be challenging to obtain a significant volume of labeled data, and it frequently takes a lot of manual labor to annotate text with the appropriate labels. According to Hugging Face (n.d.) website's inform that these pretrained models can be improved with relatively smaller quantities of labeled data because they were trained on huge datasets in the past. Because of this development, it is now more possible for developers to create natural language processing apps that work well, even with constrained datasets of resources.

## 3. Lack of Expertise

Creating natural language processing (NLP) models usually calls for a significant degree of proficiency in the fields of machine learning and natural language processing (NLP). For developers who might lack the requisite experience, this particular expertise could be a hindrance. However, pretrained models provide an answer to this problem. An example from Training transformers models in TensorFlow show that developers with little experience in natural language processing (NLP) can use these pretrained models on large datasets to create useful natural language processing applications.

**Data and Processing**

**What dataset did you use for this project?**

As in the class, most of the time we use the dataset as an excel file (CSV file), and the API keys from the website, so our group decided to use both of them to analyze which ones are more accurate on analyzing the comments from websites, and streaming platforms. For the excel file our group decided to use the IMDB movie review dataset from the website, and for the API, our group decided to use the youtube API to analyze the comments. However, if our group ask to analyze all of the comments on YouTube it is going to be hard and may cause problems with our computer so, we decide to analyze only one video by using YouTube ID from the URL which the video that we use to analyze is the music video "Santa Tell Me" from Ariana Grande.

**How did you prepare the data for analysis?**

As, our group have import the dataset from Keggel and the API that we got from website; First, we install the package that useful for our project for both CSV file and API which are pip install emoji, vander Sentiment google-api-python-client. Then, we import some of the elements from the package that have been installed which are import build (for filtering the comments), re (for filtering comments with just emojis), emoji (for analyze the sentiments of the comment), SentimentIntensityAnalyzer (for visualization), requests (for asking an access), pd (for working with data sets), KNeighborsClassifier (for learning an algorithm that makes classifications based on data neighbors), SVM (for classification or regression challenges), and lastly, accuracy_score (for analyze the accuracy). Then, request for using the API from Hugging Face that our group is going to use in this project which is LLAMA-2 from Meta. All of the processes that are

mentioned above are used for both CSV file and API Key for other procedure will be mentioned

separately below:

**CSV File:**

  After  installing the package, importing all the elements from each package, and request

from the Hugging Face already. After that, import the movie review dataset or the CSV file, the

code is `df = pd.read_csv("IMDB Dataset.csv")` which shows the outcome as shown on figure 1.

| | review | sentiment |
|---|---|---|
| 0 | One of the other reviewers has mentioned that ... | positive |
| 1 | A wonderful little production. <br /><br />The... | positive |
| 2 | I thought this was a wonderful way to spend ti... | positive |
| 3 | Basically there's a family where a little boy ... | negative |
| 4 | Petter Mattei's "Love in the Time of Money" is... | positive |
| ... | ... | ... |
| 49995 | I thought this movie did a down right good job... | positive |
| 49996 | Bad plot, bad dialogue, bad acting, idiotic di... | negative |
| 49997 | I am a Catholic taught in parochial elementary... | negative |
| 49998 | I'm going to have to disagree with the previou... | negative |
| 49999 | No one expects the Star Trek movies to be high... | negative |

50000 rows × 2 columns

*Figure 1 Outcome of the code* `df = pd.read_csv("IMDB Dataset.csv")`

  As the dataset that our group got their are massive amount of data up to fifty thousand

which this amount of data, it might cause some problem while we running other code and will

cause to our computer so, the data need to set the limit row of the data, which this will help our

project work smoothly without any problem; so the code that we use to set the limit of the row is

`limit_data = pd.read_csv("IMDB Dataset.csv", nrows=1000)` our group have set the

limitation for this dataset as one thousand rows which the outcome shown as figure 2. Then, after

setting the limit of the row which the data that have set the limit already will named as

`limit_data`, we need to to cleaning the text and count the number of emoji from the

`limit_data` but before we going to clean up and count the text, it need to set the new variable

and let the code pull out only the review row to clean up and count it which the code is `review`

`= limit_data['review']` and set the new variable as `review`. Then, start using the code on

figure 3 to clean up and count the text and emojis.

|  | review | sentiment |
|---|---|---|
| 0 | One of the other reviewers has mentioned that ... | positive |
| 1 | A wonderful little production. <br /><br />The... | positive |
| 2 | I thought this was a wonderful way to spend ti... | positive |
| 3 | Basically there's a family where a little boy ... | negative |
| 4 | Petter Mattei's "Love in the Time of Money" is... | positive |
| ... | ... | ... |
| 995 | Nothing is sacred. Just ask Ernie Fosselius. T... | positive |
| 996 | I hated it. I hate self-aware pretentious inan... | negative |
| 997 | I usually try to be professional and construct... | negative |
| 998 | If you like me is going to see this in a film ... | negative |
| 999 | This is like a zoology textbook, given that it... | negative |

1000 rows × 2 columns

*Figure 2 Outcome of the code* `limit_data = pd.read_csv("IMDB Dataset.csv", nrows=1000)`

```
hyperlink_pattern = re.compile(
    r'http[s]?://(?:[a-zA-Z]|[0-9]|[$-_@.&+]|[!*\¥(\¥),]|(?:%[0-9a-fA-F][0-9a-fA-F]))+<br>')

threshold_ratio = 0.65

relevant_comments = []

# Inside your loop that processes comments
for IMBD_review_text in review:

    IMBD_review_text = IMBD_review_text.lower().strip()

    emojis = emoji.emoji_count(IMBD_review_text)

    # Count text characters (excluding spaces)
    text_characters = len(re.sub(r'\¥s', '', IMBD_review_text))

    if (any(char.isalnum() for char in IMBD_review_text)) and not hyperlink_pattern.search(IMBD_review_text):
        if emojis == 0 or (text_characters / (text_characters + emojis)) > threshold_ratio:
            relevant_comments.append(IMBD_review_text)

# Print the relevant comments
relevant_comments[:5]
```

*Figure 3 Code for cleaning up text and count the number of text and emoji from* `limit_data`

After, clean up and count the number of text and emoji is to tokenize the IMDB movie review by using the CountVectorizer by creating a sample corpus of the dataset. Then, create the object. Next, fit the vectorizer to the corpus that has been created at first, and change the corpus into a bag of word representation. Also, create an SVM classifier and train the SVM and label as the code on figure 4 which the output shown as figure 5.

```python
#Import the CountVectorizer class
from sklearn.feature_extraction.text import CountVectorizer

from sklearn import svm
output = limit_data
output

#Create a sample corpus of documents and their labels
corpus = output['review']

labels = output['sentiment']

#Create a CountVectorizer object
vectorizer = CountVectorizer()

#Fit the vectorizer to the corpus
vectorizer.fit(corpus)

#Transform the corpus into a bag-of-words representation
bow_matrix = vectorizer.transform(corpus)

#Create an SVM classifier
model_svm = svm.SVC(C=1, kernel='linear')

#Train the classifier on the bag-of-words representation and labels
model_svm.fit(bow_matrix, labels)
```

```
▼              SVC
SVC(C=1, kernel='linear')
```

*Figure 4 Tokenization of IMDB Movie Review (SVM Classifier)*          *Figure 5 Output of the SVM Classifier*

After creating the classifier, the next step will be feature extraction and selection by using the code below as shown on figure 6 which will extract the features from the "review" column.

```python
import pandas as pd
from sklearn.feature_extraction.text import CountVectorizer

df = pd.read_csv("IMDB Dataset.csv", nrows=1000)

# Extracting features from the 'review' column using CountVectorizer
vectorizer = CountVectorizer()
X = vectorizer.fit_transform(df['review'])

# Printing the shape of the feature matrix
print(X.shape)

(1000, 17922)
```

*Figure 6 Feature Extraction and Selection*

Next, our group training and testing the data by import `train_test_split` from

`sklearn.model_selection`. Then, write the code to run the train and test function as shown on figure 7

below.

```
[ ]   from sklearn.model_selection import train_test_split
      train, test = train_test_split(df, test_size=0.2, random_state=42)
```

*Figure 7 Training and Testing Data*

Lastly, our group have imply the K-Nearest Neighbour (KNN) and Support Vector

Machine (SVM) combine with the machine learning that our group selected which LLAMA-2

which the code for running the KNN shown on figure 7, and for SVM shown figure 8

```
from sklearn.neighbors import KNeighborsClassifier
from sklearn.feature_extraction.text import CountVectorizer

X_train = train.drop('sentiment', axis=1)
y_train = train['sentiment']
X_test = test.drop('sentiment', axis=1)
y_test = test['sentiment']

knn = KNeighborsClassifier(n_neighbors = 3)
vectorizer = CountVectorizer()

X_train = vectorizer.fit_transform(train['sentiment'])
X_test = vectorizer.transform(test['sentiment'])

knn.fit(X_train, y_train)

y_pred = knn.predict(X_test)

print(y_pred[0])
y_proba = knn.predict_proba(X_test)
print(y_proba[0])

positive
[0. 1.]
```

*Figure 7 K-Nearest Neighbour Code*

```
[ ]   # Load the dataset
      df = limit_data

      # Assuming your dataset has a 'label' column for the target variable and a 'text' column for the text data
      X = df['review']
      y = df['sentiment']

      # Split the dataset into training and testing sets
      X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

      # Convert text data into numerical features using TF-IDF vectorization
      tfidf_vectorizer = TfidfVectorizer()
      X_train_tfidf = tfidf_vectorizer.fit_transform(X_train)
      X_test_tfidf = tfidf_vectorizer.transform(X_test)

      # Initialize SVM model
      svm_model = SVC(kernel='linear')

      # Train the SVM model
      svm_model.fit(X_train_tfidf, y_train)

      # Train the SVM model
      svm_model.fit(X_train_tfidf, y_train)

      # Make predictions on the test set
      y_pred = svm_model.predict(X_test_tfidf)

      # Evaluate the model
      accuracy = accuracy_score(y_test, y_pred)
      classification_report_result = classification_report(y_test, y_pred)

      # Print the results
      print(f"Accuracy: [accuracy]")
      print("Classification Report:")
      print(classification_report_result)
```

```
Accuracy: 0.815
Classification Report:
              precision    recall  f1-score   support

    negative       0.82      0.83      0.82       104
    positive       0.81      0.80      0.81        96

    accuracy                           0.81       200
   macro avg       0.81      0.81      0.81       200
weighted avg       0.81      0.81      0.81       200
```

*Figure 8 Support Vector Machine Code*

**API Key:**

First, after install all the package that need to use on this project and request for using the LLAMA-2 from Hugging Face, our group import the API key for YouTube API, and copy the URL from youtube any video that you want to analyze the comment which our group chose "Santa Tell Me" from Ariana Grande as a sample for this project as the code below on figure 9.

```
[ ]  API_KEY = "AIzaSyD_n2OzFTMPtxXLaVAn5Zj-_I3vIdEAU4M"# Put in your API Key

     youtube = build('youtube', 'v3', developerKey=API_KEY) # initializing Youtube API

     # Taking input from the user and slicing for video id
     video_id = "nIROMkrRklg"
     # input('Enter Youtube Video URL: ')[-11:]
     # https://www.youtube.com/watch?v=nIROMkrRklg
     print("video id: " + video_id)

     # Getting the channelId of the video uploader
     video_response = youtube.videos().list(
       part='snippet',
       id=video_id
     ).execute()

     # Splitting the response for channelID
     video_snippet = video_response['items'][0]['snippet']
     uploader_channel_id = video_snippet['channelId']
     print("channel id: " + uploader_channel_id)

     video id: nIROMkrRklg
     channel id: UCOVOyT2OCBKdQhF3BAbZ-1g
```

*Figure 9 YouTube API code*

As the comments on the selected video have over thirteen thousand comments so, it needs to be minimized to avoid the problem that will cause to our project so our group decided to set the limit for the number of comments as one thousand as a code below (Figure 10), an example of output on figure 11 and analyze the items. Also, input the code to analyze each dataset column, the non-null count and also the data type object as it is shown on figure 12.

*Figure 10 Limit the YouTube Comments*



*Figure 11 Outcome of the Limit Comment Code*



*Figure 12 Analyze Dataset Column Code*

After analyze each column dataset, our group clean up text and count the number of the emoji same code as the IMDB movie review data python code on figure 3 but the different is change the variable of the code which the code are shown below (Figure 13).

```
[ ]  hyperlink_pattern = re.compile(
       r'http[s]?://(?:[a-zA-Z]|[0-9]|[$-_@.&+]|[!*¥¥(¥¥),]|(?:%[0-9a-fA-F][0-9a-fA-F]))+<br>')

     threshold_ratio = 0.65

     relevant_comments = []

     # Inside your loop that processes comments
     for comment_text in comments:

       comment_text = comment_text.lower().strip()

       emojis = emoji.emoji_count(comment_text)

       # Count text characters (excluding spaces)
       text_characters = len(re.sub(r'¥s', '', comment_text))

       if (any(char.isalnum() for char in comment_text)) and not hyperlink_pattern.search(comment_text):
         if emojis == 0 or (text_characters / (text_characters + emojis)) > threshold_ratio:
           relevant_comments.append(comment_text)

     # Print the relevant comments
     relevant_comments[:5]
```

*Figure 13 Code for cleaning up text and count the number of text and emoji for YouTube API*

Next, analyze each comment on the video by using `sentiment_scores` and set the sentiment score as figure 14; for the score, it can adjust in any number as the owner wants, which our group decides to set the sentiment score as 0.05. If the comment has been calculated and the result is more than 0.05 that comment will be a positive commnet. However, if the score is less than 0.05 that comment will be Negative comment automatically, and if the comment is equal to 0.05 the comment will be neutral which will be the outcome as shown on figure 15.

```python
def sentiment_scores(comment):

    # Creating a SentimentIntensityAnalyzer object.
    sentiment_object = SentimentIntensityAnalyzer()
    sentiment_dict = sentiment_object.polarity_scores(comment)

    return sentiment_dict['compound']

polarity = []
positive_comments = []
negative_comments = []
neutral_comments = []

comments = relevant_comments
print("Analysing Comments...")
for index, items in enumerate(comments):
    # Append the compound polarity score of the comment to the polarity list.
    polarity.append(sentiment_scores(items))

    # Determine the sentiment of the comment.
    if float(polarity[-1]) > 0.05:
        sentiment = "Positive"
        positive_comments.append(items)
    elif float(polarity[-1]) < -0.05:
        sentiment = "Negative"
        negative_comments.append(items)
    else:
        sentiment = "Neutral"
        neutral_comments.append(items)

    # Assign the sentiment to the sentiment column of the df DataFrame.
    df.loc[index, 'sentiment'] = sentiment

# Print polarity
print(polarity[:5])

# Print sentiment column
print(df['sentiment'])
```

*Figure 14 Sentiment Score Code*

```
Analysing Comments...
[0.4404, 0.4404, 0.0, 0.5423, 0.0]
0      Positive
1      Positive
2       Neutral
3      Positive
4       Neutral
         ...
88     Positive
89     Positive
90     Negative
91      Neutral
92          NaN
Name: sentiment, Length: 93, dtype: object
```

*Figure 15 Outcome of The Sentiment Code*

After running the sentiment score, our group used the K-Nearest Neighbour (KNN) and Support Vector Machine (SVM) combined with the machine learning that had been requested at first. Which the KNN and SVM that our group is to calculate the accuracy of the sentiment score. But before running both KNN model and SVM model, it needs testing and training data first as shown on figure 16.

```
from sklearn.model_selection import train_test_split
train, test = train_test_split( comment_text, test_size = 0.33, random_state=42)
```

*Figure 16 Training and Testing Data Code*

Then, after testing and training, the KNN model and SVM model are ready to use to calculate the accuracy. For the KNN model is going to use the code that shown on figure 17, and for the SVM model will use the code as shown on figure 18, which the outcome of KNN and SVM model are different; for KNN the output shown as `Accuracy: 0.0` and for SVM model the outcome is `Accuracy: 1.0`

As the result of accuracy check are not the same so our group think that the reason is when the code run the sentiment score are not reliable enough because the sentiment score can be change depend on users which out decide to use 0,05 but however it can be adjust to any number that users want to which made the code of KNN and SVM show different output as it explain above.

## Model Selection and Training

Hugging Face contains tons of machine learning models as an API Key. So our group chose LLAMA-2. As Hugging Face mentioned, LLAMA-2 is "a collection of pretrained and fine-tuned generative text models ranging in scale from 7 billion to 70 billion parameters. This is the repository for the 7B fine-tuned model, optimized for dialogue use cases and converted for the Hugging Face Transformers format" (Hugging Face, n.d.). Links to other models can be found in the index at the bottom.There are two main reasons that our group decided to use this model: first, as we look through the download counts and like counts, there are a lot of people who use this model for their work or project, so our group think that this model is good to work with our project, and reliable. Second reason is LLAMA-2 have been create and develop by Meta that release this model for public users which LLAMA-2 is Llama 2 family of large language models (LLMs) also, this model contain fine tune generative text model range between seven billion to seventy billion.

The processes that our group uses to train this model are K-Nearest Neighbour (KNN) and Support Vector Model (SVM) to check the accuracy of the code; rather the code analyzes the comments accurately or not.

**Results and Discussion**

**What did you learn from your analysis?**

Overall, the things that our group gain from this project basically this project let us know that the machine learning not work with only the excel file only but it can also use with API Keys which both them can be imply with the model that our group chose. However, as our group doing the coding it also let us know that analyze the comment or the review better use the machine learning model with Excel or the CSV will be accurate than using API Key because the API Key we need to set the sentiment score by ourselves which make the users able to set their own score and let the code run which have low accuracy than using CSV file base on our group work experience.

**What are the implications of your findings?**

Our group decided to use the machine learning model with review and comments analytic which will help to separate good comments and bad comments (ex. Hate speech, bad word). This comment analyzes code able to be used with any website or webpage that allow public users to share their own thoughts such as Instagram, Facebook, YouTube, and TikTok. This kind of comment analytic are able to develop and use to detect those bad or negative comments on any website that allow people to share their thoughts through social media.

**Conclusion and Future Work**

**What are the main takeaways from your research?**

The main takeaways from our research from the CSV file and API Key need to chose the dataset that have enough data to analyze which rarely see on CSV file, but for the YouTube API Key, it need to find the video that have enough comments to analyze if the the video have less amount of comment the code will not able to run it properly base on our group experience. Furthermore, when the machine learning model test and train with KNN and SVM model the CSV file will show the accuracy proper than the API Key that the user need to set their own sentiment score by themselves.

**What are some potential future directions for research using Hugging Face pretrained models?**

Improving the robustness and generalization of pretrained models is crucial for their successful application across various domains, languages, and tasks. For instance, adapting s model initially trained on English news articles to generate summaries for Japanese medical reports requires a nuanced approach. This might involve additional training on domain-specific datasets to understand medical terminology and context. Similarly, fine-tuning a model for low-resource languages and dialects, and incrementally training the model to capture the linguistic nuances unique to that language. Such strategies are essential for extending the capabilities of natural language processing models beyond their initial training scope, enabling them to provide valuable insights and perform tasks in new and diverse settings.

Furthermore on another topic such as, enhancing the creativity and expressiveness of pretrained models is a fascinating area of development in natural language processing. It involves training models to generate text that is not only novel and engaging but also varied in

style, tone, and emotion. This could involve generating text that continues a story based on the reader's choices or adapting the formality of language to suit the user's communication style. The key to these advancements is a combination of innovative training techniques, diverse datasets, and iterative fine-tuning to capture the subtleties of human expression and interaction.

**References**

Maning, J. (2023, September 22). *What Is Hugging Face and What Is It Used For?*. Retrieved

from https://www.makeuseof.com/what-is-hugging-face-and-what-is-it-used-for/

Nallapati, R., Zhou, B., dos Santos, C. N. d., Gulcehre, C., & Xiang, B. (2016). *Abstractive Text*

*Summarization Using Sequence-to-Sequence RNNs and Beyond.* Retrieved from

http://arxiv.org/abs/1602.06023

Fu, Z., Yang, H., So, A. M.-C., Lam, W., Bing, L., & Collier, N. (2022). *On the Effectiveness of*

*Parameter-Efficient Fine-Tuning.* Retrieved from https://arxiv.org/abs/2211.15583

Wolf, T., Debut, V., Sanh, V., Chaumart, J., & Rush, A. M. (2019). *Hugging Face's*

*Transformers: State-of-the-art Natural Language Processing*. Retrieved from

https://arxiv.org/abs/1910.03771

Devlin, J., Chang, M.-W., Lee, K., Toutanova, K., & Lee, J. (2018). BERT: Pre-training of Deep

Bidirectional Transformers for Language Understanding. Retrieved from

https://arxiv.org/abs/1810.04805

Han, X., Zhang, Z., Ding, N., Gu, Y., Liu, X., Huo, Y., Qiu, J., Yao, Y., Zhang, A., Zhang, L.,

Han, W., Huang, M., Jin, Q., Lan, Y., Liu, Y., Liu, Z., Lu, Z., Qiu, X., Song, R., Tang, J.,

Wen, J.-R., Yuan, J., Zhao, W. X., & Zhu, J. (2021). *Pre-Trained Models: Past, Present*

*and Future.* https://arxiv.org/abs/2106.07139

Hugging Face. (n,d,). Llama-2-7b-chat-hf. Hugging Face.

https://huggingface.co/meta-llama/Llama-2-7b-chat-hf

Hugging Face. (n.d.). *Pretrained models*. Retrieved April 5, 2023, from

https://huggingface.co/transformers/v3.3.1/pretrained_models.html

Training transformers models in TensorFlow. (n.d.). Retrieved April 5, 2023, from

https://colab.research.google.com/github/huggingface/notebooks/blob/main/transformers

_doc/en/tensorflow/training.ipynb