

✓ Install Required Libraries:

```
1 # @title Install Required Libraries:
2 !pip install transformers scikit-learn PyMuPDF
```

```
Requirement already satisfied: transformers in /usr/local/lib/python3.10/dist-packages (4.46.3)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.10/dist-packages (1.5.2)
Requirement already satisfied: PyMuPDF in /usr/local/lib/python3.10/dist-packages (1.25.1)
Requirement already satisfied: filelock in /usr/local/lib/python3.10/dist-packages (from transformers) (3.16.1)
Requirement already satisfied: huggingface-hub<1.0,>=0.23.2 in /usr/local/lib/python3.10/dist-packages (from transformers) (0.26.5)
Requirement already satisfied: numpy>=1.17 in /usr/local/lib/python3.10/dist-packages (from transformers) (1.26.4)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from transformers) (24.2)
Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.10/dist-packages (from transformers) (6.0.2)
Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.10/dist-packages (from transformers) (2024.9.11)
Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (from transformers) (2.32.3)
Requirement already satisfied: tokenizers<0.21,>=0.20 in /usr/local/lib/python3.10/dist-packages (from transformers) (0.20.3)
Requirement already satisfied: safetensors>=0.4.1 in /usr/local/lib/python3.10/dist-packages (from transformers) (0.4.5)
Requirement already satisfied: tqdm>=4.27 in /usr/local/lib/python3.10/dist-packages (from transformers) (4.66.6)
Requirement already satisfied: scipy>=1.6.0 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (1.13.1)
Requirement already satisfied: joblib>=1.2.0 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (1.4.2)
Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (3.5.0)
Requirement already satisfied: fsspec>=2023.5.0 in /usr/local/lib/python3.10/dist-packages (from huggingface-hub<1.0,>=0.23.2->transformers) (2024.9.11)
Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.10/dist-packages (from huggingface-hub<1.0,>=0.23.2->transformers) (4.12.2)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests->transformers) (3.3.2)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests->transformers) (3.10)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests->transformers) (2.2.3)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests->transformers) (2024.8.30)
```

✓ Extract Text from PDFs:

```
1 # @title Extract Text from PDFs:
2 import fitz # PyMuPDF
3
4 def extract_text_from_pdf(pdf_path):
5     doc = fitz.open(pdf_path)
6     text = ""
7     for page_num in range(len(doc)):
8         page = doc.load_page(page_num)
9         text += page.get_text()
10    return text
11
12 # Example PDF paths
13 pdf_paths = ["butterflies.pdf", "docker.pdf", "flowers.pdf", "rockmusic.pdf"]
14
15 # Extract text from each PDF
16 texts = [extract_text_from_pdf(pdf_path) for pdf_path in pdf_paths]
```

✓ Generate Embeddings:

```
1 # @title Generate Embeddings:
2 from transformers import AutoTokenizer, AutoModel
3 import torch
4
5 # Load pre-trained model and tokenizer
6 model_name = 'sentence-transformers/all-MiniLM-L6-v2'
7 tokenizer = AutoTokenizer.from_pretrained(model_name)
8 model = AutoModel.from_pretrained(model_name)
9
10 # Function to generate embeddings
11 def get_embeddings(texts):
12     inputs = tokenizer(texts, padding=True, truncation=True, return_tensors='pt')
13     with torch.no_grad():
14         outputs = model(**inputs)
15     return outputs.last_hidden_state.mean(dim=1).numpy()
16
17 # Generate embeddings for the extracted texts
18 embeddings = get_embeddings(texts)
```

```

/usr/local/lib/python3.10/dist-packages/huggingface_hub/utils/_auth.py:94: UserWarning:
The secret `HF_TOKEN` does not exist in your Colab secrets.
To authenticate with the Hugging Face Hub, create a token in your settings tab (https://huggingface.co/settings/tokens), set it as :
You will be able to reuse this secret in all of your notebooks.
Please note that authentication is recommended but still optional to access public models or datasets.
warnings.warn(
```

✓ Cluster Embeddings:

```

1 # @title Cluster Embeddings:
2 from sklearn.cluster import KMeans
3
4 # Number of clusters
5 num_clusters = 3
6
7 # Apply KMeans clustering
8 kmeans = KMeans(n_clusters=num_clusters, random_state=0)
9 clusters = kmeans.fit_predict(embeddings)
10
11 # Print cluster assignments
12 for pdf_path, cluster in zip(pdf_paths, clusters):
13     print(f"PDF: {pdf_path} -> Cluster: {cluster}")

```

```

➡ PDF: butterflies.pdf -> Cluster: 0
   PDF: docker.pdf -> Cluster: 2
   PDF: flowers.pdf -> Cluster: 0
   PDF: rockmusic.pdf -> Cluster: 1

```

✓ Prepare Data for Visualization: Save the clustering results to a JSON file that can be used by the D3.js script.

```

1 # @title Prepare Data for Visualization: Save the clustering results to a JSON file that can be used by the D3.js script.
2 import json
3
4 # Prepare data for visualization
5 data = [{"pdf": pdf_path, "cluster": int(cluster)} for pdf_path, cluster in zip(pdf_paths, clusters)]
6
7 # Save data to a JSON file
8 with open("clustering_results.json", "w") as f:
9     json.dump(data, f)

```



```

1 import numpy as np
2 import plotly.express as px
3 import pandas as pd
4
5 # Prepare data for visualization
6 data = [{"pdf": pdf_path, "cluster": int(cluster)} for pdf_path, cluster in zip(pdf_paths, clusters)]
7
8 # Generate random coordinates for demonstration purposes
9 for d in data:
10     d['x'] = np.random.rand()
11     d['y'] = np.random.rand()
12     d['z'] = np.random.rand()
13
14 # Convert data to a DataFrame
15 df = pd.DataFrame(data)
16
17 # Create a 3D scatter plot
18 fig = px.scatter_3d(df, x='x', y='y', z='z', color='cluster', hover_data=['pdf'], title='PDF Clustering Visualization')
19
20 # Add data labels
21 fig.update_traces(marker=dict(size=5),
22                   selector=dict(mode='markers'))
23
24 # Show the plot
25 fig.show()

```



PDF Clustering Visualization

