# co cu # ta # # # # # # # # # # # # # # # # #	Step 1 - Reading the Tables from Database file  mport sqlite3  Connect to the Sqlite database onn = sqlite3.connect('eng.subtitles_database.db') ursor = conn.cursor()  Execute the query to fetch table names ursor.execute("SELECT name FROM sqlite_master WHERE type='table'")  Fetch all table names from the database able_names = cursor.fetchall()  Print the table names intrit(table names)
[(': S 4]: cu co fo num name con:	
[5]: df df [5]:  0 1 2 3	### ### ##############################
<pre><cla #="" 0="" 1="" 2="" data="" dty ="" memore<="" pre="" rang=""></cla></pre>	num 82498 non-null int64 name 82498 non-null object content 82498 non-null object vpes: int64(1), object(2) nory usage: 1.9+ MB  Unzipping the content of row and decoding using latin-1  mport random
im im  # de	Function to decode binary data from ZIP archives ef decode_method(binary_data):     with io.BytesIO(binary_data) as f:         with zipfile.Zipfile(f, 'r') as zip_file:
df # fo  # im im # #	Create a new column to store the selected rows  f['file content'] = ''  Iterate over selected rows and populate the "file content" column or idx in selected_rows_indices:
un. # nul # se	<pre>with zipfile.ZipFile(f, 'r') as zip_file:     # Assuming there's only one file in the ZIP archive     subtitle_content = zip_file.read(zip_file.namelist()[0]) # Now 'subtitle_content' should contain the extracted subtitle content return subtitle_content.decode('latin-1') # Assuming the content is UTF-8 encoded text  Get unique file names from the DataFrame nique_file_names = df['name'].unique().tolist()  Calculate the number of files to select um_files_to_select = int(0.3 * len(unique_file_names))  Randomly select files without replacement elected_file_names = random.sample(unique_file_names, num_files_to_select)  Create a new column to store the selected files' content f['file content'] = ''</pre>
fo	Iterate over selected files, decode their content, and populate the "file content" column or file_name in selected_file_names: file_rows = df[df['name'] = file_name] for _, row in file_rows.iterrows():     decoded_content = decode_method(row['content'])     df.at[row.name, 'file content'] = decoded_content  Now df contains the original data along with the decoded content for approximately 30% of the files  f  num
82 82 82	2 9180592 yumis.cells.s02.e13.episode.2.13.(2022).eng.1cd b'PK\x03\x04\x14\x00\x00\x00\x00\x08\x00L\xb9\x99V 3 9180592 yumis.cells.s02.e14.episode.2.14.(2022).eng.1cd b'PK\x03\x04\x14\x00\x00\x00\x00\x00\x08\x00L\xb9\x99V 4 9180600 broker.(2022).eng.1cd b'PK\x03\x04\x14\x00\x00\x00\x00\x00\x00\x08\x00L\xb9\x99V 5 2493 521935 the.prophets.game.(2000).eng.1cd b'PK\x03\x04\x14\x00\x00\x00\x00\x00\x00\x00\x00\x00\x0
[9]: im im im im	### Assuming there's only one file in the ZIP archive file_content_decode('latin-1') # Assuming the content is UTF-8 encoded text exception as e:
	print("Error:", e) return None  Define a function to randomly select 30% of files and extract their content ef randomly_select_and_extract(df):     # Get unique_file_names from the DataFrame     unique_file_names = df['name'].unique().tolist()  # Calculate the number of files to select     num_files_to_select = int(0.3 * len(unique_file_names))  # Randomly select files without replacement     selected_file_names = random.sample(unique_file_names, num_files_to_select)  # Create a new column to store the selected files' content df['file_content'] = ''  # Iterate over selected files, decode their content, and populate the "file content" column
df # ,df.	for file_name in selected_file_names:     file_rows = df[df['name'] == file_name]     for _, row in file_rows.iterrows():         decoded_content = extract_zip_content(row['content'])         df.at[row.name, 'file_content'] = decoded_content  return df  Read the data into a pandas DataFrame (replace "conn" with your database connection)     f = pd.read_sql_query("""SELECT * FROM zipfiles""", conn)  Apply the random selection and extraction function     f_with_random_selection = randomly_select_and_extract(df)  Now df_with_random_selection contains the original data along with the decoded content for approximately 36% of the files  f_with_random_selection
	Num
82 82 824	### PK\x83\x84\x14\x80\x80\x80\x80\x80\x80\x80\x80\x90\x99\x99\  ################################
82 82 82 82 82 Na	b'PK\x83\x84\x14\x89\x89\x88\x88\x88\x  2493 b'PK\x83\x84\x14\x89\x89\x88\x88\x88\x  2494 b'PK\x83\x84\x14\x89\x89\x88\x88\x88\x  2495 b'PK\x83\x84\x14\x89\x89\x88\x88\x88\x88\x88\x  2497 b'PK\x83\x84\x14\x89\x89\x89\x88\x88\x88\x88\x88\x88\x  2497 b'PK\x83\x84\x14\x89\x89\x89\x88\x88\x88\x88\x88\x88\x88
df. # pr  0 9 1 9 3 9 4 9	Pecode the content in the 'content' column and create a new column for the decoded content  f_with_random_selection['decoded_content'] = df_with_random_selection['content'].apply(decode_content)  Display the DataFrame with the decoded content column  rint(df_with_random_selection.head())  num  name \ 9180533
2   3   4   1   1   2   1   3   1   3   1   3   1   3   1   3   1   1	b'PK\x83\x04\x14\x00\x08\x00\x08\x00\x09\x99\ b'PK\x83\x04\x14\x00\x00\x00\x00\x00\x09\x99\ b'PK\x03\x04\x14\x00\x00\x00\x00\x00\x00\x00\x00\x00\x0
L3]:	Num
82 82 82 824 14]: df.	### 15/21935   ### 16-prophets game (2000).eng.1cd   bPK\x03\x04\x14\x00\x00\x00\x00\x08\x00\x00\x08\x00\x00
82 DZ 82 82 0 0 82 0 < 82 Na	PKIDDOU@DVYRD@YYD06Yumi's.Cells.S02   PKIDDOU@DVYRD@YYD06Yumi's.Cells.S02   PKIDDOU@DVYRD@YYD06Yumi's.Cells.S02   PKIDDOU@DVYRD@YYD06Yumi's.Cells.S02   PKIDDOU@DVYRD@YYD06Yumi's.Cells.S02   Call the properties of the properties o
df. # df.	decoded_content = content.decode('latin-1').encode('utf-8', errors='ignore').decode('utf-8')     return decoded_content     except Exception as e:         print("Error decoding content:", e)     return None  Decode the content in the 'content' column and create a new column for the decoded content  f_with_random_selection['decoded_content'] = df_with_random_selection['content'].apply(decode_content)  Convert binary data in 'decoded_content' column to readable text  f_with_random_selection['decoded_content'] = df_with_random_selection['decoded_content'].apply(lambda x: x if isinstance(x, str) else '')  Display the DataFrame with the decoded content column  rint(df_with_random_selection['decoded_content'].head())  PKUBLUB VIXID &BUBLUS (The. Message.1976  PKUBLUB VIXID &BUBLUS (The. Message.1976  PKUBLUB VIXID &BUBLUS (The. Message.1976)
# de	PMIDUUMEN'RNUAYYOUR'SUN'S.Cells.Se2 PMIDUUMEN'PROKEY 2022 720p  me: decoded_content, dtype: object  Define a function to format the decoded content with subtitle numbers  ef format_decoded_content(row):     subtitle_number = row['num']     decoded_content = row['decoded_content']     formatted_content = f"Subtitle {subtitle_number} ({row['name']}):\n\n{decoded_content}"     return formatted_content  Apply the formatting function to create a new column with formatted content f_with_random_selection['formatted_content'] = df_with_random_selection.apply(format_decoded_content, axis=1)  Display the DataFrame with the formatted content rint(df_with_random_selection['formatted_content'].head())  Subtitle 9180533 (the.message.(1976).eng.icd): Subtitle 9180533 (there.comes.the.grump.s01.e09
16]: df. 16]: 0 1 2 3 4 82 82 82	Subtitle 9180592 (yumis.cells.s02.e13.episode Subtitle 9180594 (yumis.cells.s02.e14.episode Subtitle 9180600 (broker.(2022).eng.1cd):\n\nP ne: formatted_content, dtype: object  f_with_random_selection['formatted_content']
[17]: im	Extractind data which present inside the content column  mport re  Function to extract content inside subtitles ef extract_content_inside(formatted_content):  # Define a pattern to extract content inside the subtitles pattern = r'\d+\n(.'7)(?=\d+\n[s])*  # Use regular expression to find matches matches = re.findall(pattern, formatted_content, flags=re.DOTALL)  # Join the matches with newline separator content_inside = '\n'.join(matches) return content_inside
df. # pr  0 1 1 2 3 4	Apply the function to extract content inside each row [_with_random_selection['content_inside'] = df_with_random_selection['formatted_content'].apply(extract_content_inside)  Display the DataFrame with the newly added content_inside column rint(df_with_random_selection[['name', 'content_inside']].head())  name
18]: 0 0ç Øl 1 2 ¢X 3 4 vD A	+680 vD+±UYC> ÑUA eQ`d)ÂO>-ĒFU\rx"G%U0>O, ÂU^  f_with_random_selection['content_inside']  0 cO0 îŭur+U0 ¶U vî;s'cóweoôO,U y.U3 z U ofAU SYOO8 a  XOU 3U ±80 dU Oño ZdUU A±4rjiĒĒEX£āOcu U x56 <cmu +60="" -u="" ?="" bj="U" b}v="" d+±uyc="" içou="" mi="" u="" wfu="" xu="" x¶nu="" äääou="tu" āu="" āñu=""> ÑU D+±UYC&gt; Ñ</cmu>
82 82 82 0 Î Na 13]: im im # /	Trill_rid_Cqn_=0ell_1All_All_All_All_All_All_All_All_All_
0 1 2 3 4	print("Error decoding content:", e) return None  Apply the decode_content function to decode the 'content_inside' column f_with_random_selection['decoded_content_inside'] = df_with_random_selection['content_inside'].apply(decode_content)  Display the DataFrame with the newly added 'decoded_content_inside' column rint(df_with_random_selection[['name', 'decoded_content_inside'].head())  Aname
3 4 4 24]: df. 24]: 0 0ç ø[ 1 2 ¢X 3 4 VD A	AIT exôt 30 ± ±0 dt 0 the Zdtt A±drjiÉEExRaôcul Isse <cmu *="" +="" -="" 0="" 1="" 3)v="" 50="" 680="" aadd="" ait="" c="" etd="" iáhn="" lçobu="" mytu="" ut="" v="" wft="" wi="" xo="" ±=""> [N A e Q' d) Áo &gt; -ÉFO \ r x * ~ 6 w U &gt; 0, du ~  f_with_random_selection['decoded_content_inside']  côt îver * Ut 1 the Young Ait of the Young Ait of</cmu>
18 82 82 82 82 82 19 1 im	az "DaenĀÑSIK n̄ aκf <a "in="" "oub!"="" 'nāqshi="" **="" **<="" *wo="" 1,="" 10="" 2494="" 2495="" 2496="" 65="" \rww.xióe\r!id2qen="óed" aud="" dudf="" extrictor="" ko="" out="" td="" uikiturāī="" y8āqā'="" ~¿ée.00="" ºðo="" āilā="" ōoxemþ±¿10="" δ\0="" ἀnμαu,=""></a>
ba of row who	<pre>Query the 'zipfiles' table to fetch rows in batches atch_size = 1000 ffset = 0 pws = [] hile True:     cursor.execute(f'SELECT * FROM zipfiles LIMIT {batch_size} OFFSET {offset}')     batch_rows = cursor.fetchall()     if not batch_rows)         break     rows.extend(batch_rows)     offset += batch_size  Close the cursor and connection to the database     ursor.close()     ponn.close()     Calculate 30% of the total number of rows     urrows.to_process = int(len(rows) * 0.3)</pre>
# da de # de	Initialize an empty list to store the data  ata = []  Define a function to extract dialogues from content  ef extract_dialogues(content):     dialogue_pattern = r' (\d(2)\d(2)\d(2)\d(2)\d(2)\d(2)\d(2)\d(2
df	<pre># Assuming there is only one file in each ZIP archive extracted_file = zip_ref.open(zip_ref.namelist()[6]) # Read and decoded the content of the extracted file with fallback encoding try:</pre>
0 1 2 3 4  2474 2474 2474 2474 2474	9279492 scorpion.s04.e15.wave.goodbye.(2018).eng.1cd 745 9279493 scorpion.s04.e16.nerd.wind.fire.(2018).eng.1cd 746 9279494 scorpion.s04.e17.dumbster.fire.(2018).eng.1cd 747 9279495 scorpion.s04.e18.dork.day.afternoon.(2018).eng
2474 2474 2474 2474 [24 18]: im im im im im	00:00:06,000> 00:00:12,074\nWatch any video  144 00:00:01,101> 00:00:03,798\n- WALTER: <i>15PE  145 00:00:01,101&gt; 00:00:04,367\n- WALTER:<i 00:00:01,101="" 146="" 15pe=""> 00:00:04,367\n- WALTER:<i 00:00:02,002="" 147="" 15pe=""> 00:00:04,071\n- WALTER:<i 00:00:01,101="" 148="" 15pe=""> 00:00:04,071\n- WALTER:<i 00:00:02,002="" 149="" 15pe=""> 00:00:05,029\n- Kid, we'll ru  140 00:00:01,101&gt; 00:00:02,399\nWALTER:<i 00:00:02,002="" 141="" 15pe=""> 00:00:02,399\nWALTER:<i ('eng_subtitles_database.db')<="" 142="" 143="" 144="" 145="" 146="" 147="" 15pe="" 3="" columns]="" database="" in="" nport="" rows="" squite="" td="" the="" to="" x="" zipfile=""></i></i></i></i></i></i></i>
# cu # rou cu co # #	<pre>query the 'zipfiles' table to fetch all rows ursor.execute('SELECT * FROM zipfiles')  Fetch all rows from the table ows = cursor.fetchall()  Close the cursor and connection to the database ursor.close() onn.close()  Define a function to extract dialogues from content ef extract_dialogues(content):     dialogue_pattern = r'\d(2)\d(2)\d(2)\d(2)\d(2)\d(2)\d(2)\d(2)</pre>
da #	return extracted_content  Initialize an empty list to store the data ata = []  Example of processing subtitle files or row in rows: num, name, content = row # Read the ZIP file content from the database zin_file = io.BytesIO(content) # Extract the contents of the ZIP file with zipfile zipfile(zip_file, 'r') as zip_ref: # Assuming there is only one file in each ZIP archive extracted_file = zip_ref.open(zip_ref.namelist()[0])  try: # Read and decode the content of the extracted file decoded_content = extracted_file.read().decode('utf-8') except UnicodeDecodeError:
	except UnicodeDecodeError:  try:  # Try decoding with a different encoding decoded_content = extracted_file.read().decode('latin1') except UnicodeDecodeError:  # Handle the decoding error gracefully print(f"Error decoding file: {name}. Skipping") continue  # Extract dialogues from the decoded content separated_content = extract_dialogues(decoded_content) # Append data to the list data.append((num, name, separated_content))  Convert the list of tuples to a DataFrame f = pd.DataFrame(data, columns=['Subtitle Number', 'Name', 'Separated Content'])  Print the DataFrame init(df)
df	Subtitle Number Name \
df	Subtitle Number 9180533 the.message.(1976).eng.1cd 9180533 here.comes.the.grump.s01.e09.joltin.jack.in.bo 9180592 yumis.cells.s02.e13.episode.2.13.(2022).eng.1cd 9180594 yumis.cells.s02.e14.episode.2.14.(2022).eng.1cd 9180600 broker.(2022).eng.1cd 9180600 broker.(2022).eng.1cd 9180600 broker.(2022).eng.1cd 9521935 the.prophets.game.(2000).eng.1cd 9521937 west.beirut.(1998).eng.1cd 9521938 frankenstein.the.true.story.(1973).eng.1cd 9521938 frankenstein.the.true.story.(1973).eng.1cd 9521940 frankenstein.the.true.story.(1973).eng.1cd
0 1 2 3 4 8248 8248 8248 8248 8248 8248 824	Subtille Number
0 1 2 3 4 824 824 824 824 824 824 824 824 824 824	Subtilic Number

