### Big Data Analytics 2024 Project 1

:movie_camera: Welcome to **CineSense**

**CineSense** is an innovative video-processing startup that extracts valuable insights from social media video content. It uses advanced natural language processing (NLP) and computer vision techniques to analyse videos' sentiments and emotions. This analysis is crucial for businesses seeking to understand their audience, improve customer experiences, and make data-driven decisions.

On your first day at **CineSense** as a data engineer, you were assigned a critical project: **find the best strategy to download and analyse videos**. **CineSense** uses YouTube videos posted by its target audience. Each video must be analysed, and results should be extracted and shared with customers. Usual tasks include extracting video, audio, transcripts, sentiments, and emotions. Since you are paired with a customer from Madrid, you will also need to transcribe videos in Spanish.

Thanks to Python, you can showcase your skills and automate the tasks. Your project leader is particularly looking for a strategy to complete tasks faster! Parallel processing concepts can help you speed up your analysis tasks, contributing to your company's success. Your colleagues prepared a project description (**part 1**) and a set of skeleton scripts for your developments provided (**part 2**).

Good luck :rocket:

### Part 1: Project description

* Develop a Python application using multiprocessing, threading or asynchronous programming concepts to download and analyse YouTube videos.

17 * Compare and contrast different solutions for serial and parallel processing executions.

18

19 * You have four weeks to complete all the tasks and document your solutions.

20

21 ##### Tasks:

22

23 **1.** Manually retrieve 10-15 random video URLs from YouTube.

24

25     * Save the URLs in a text file called `video_urls.txt`, where each URL should be stored on a separate line.

26     * Consider YouTube videos that are 2-3 minutes in duration.

27

28 **2.** Develop a Python script to read the URLs.

29

30     - Assuming you have the text file named `video_urls.txt` containing the URLs of YouTube videos, load it in Python and extract the URLs using your preferred data structure.

31

32     **[2.5 marks]**

33

34 **3.** Develop a Python script to download the videos using their URLs.

35

36     - Test your solution by downloading the files serially.

37     - Use parallel programming such as multiprocessing or threading to handle downloads. Your decision will determine the best strategy.

38     - For testing reasons, ensure the script can download up to `5` videos simultaneously to avoid YouTube blocks.

39     - You are advised to use `threads` and `semaphores` to control the downloads.

40     - Compare serial and parallel executions for your video download script.

```
41        - Discuss the complexity of your video download
   scripts' time and space.

42

43    **[15 marks]**

44

45 4. Develop a Python script to keep a log for each
   download.

46

47    - After downloading each video, create a logger to
    record which video was downloaded by which process
   or thread.

48

49    - Save the log entries to the same file, e.g., `
   download_log.txt`.

50

51    - For this script, you have to use `threads` and a
    `mutex`.

52

53    - The entries could be in the following format:

54

55    - ```json

56        "Timestamp": 12:23, 21 May 2024, "URL":"http
   ://www.youtube.com/1234", "Download":True

57        "Timestamp": 12:25, 21 May 2024, "URL":"http
   ://www.youtube.com/1235", "Download":True

58        ```

59

60

61    **[12.5 marks]**

62

63 5. Develop Python scripts to perform various video
   analysis tasks.

64

65    - After downloading a video, perform the following
    tasks.

66    - It is preferable to develop a separate script
   for each functionality.

67    - The five analysis subtasks that you have to
   develop include the following:

68        - Extract audio from a video file.

69        - Transcribe audio to text.
```

70        - Perform the sentiment analysis on a video's content, extracting its polarity and sensitivity.
71        - Translate the text into another language, e.g. Spanish.
72        - Extract the emotions of a text.
73     - Each output task should store its results in a dedicated folder designated for each video, using the video title. Feel free to organise your folder structure as you prefer.
74     - You can use any library, including `moviepy` for loading video and `speech_recognition` or `textblob` for sentiment analysis.
75     - To implement the analysis subtasks, you must use at least one of the following libraries: `multiprocessing`, `threading`, or `asyncio`.
76     - You must compare serial, multiprocessing, threading, and concurrency for at least one of the subtasks, such as the extracting audio functionality. You do not have to do it for the rest of the subtasks.
77
78     **[30 marks]**
79
80 **6.** Document your solutions in a short technical report.
81
82     * This project is open to your interpretations, and you can use any programming strategy that fits better to this project, including processes, threads, or asynchronous programming.
83     * Discuss briefly your strategy for developing each solution.
84     * A 1-2 page discussion (600-1200 words) will suffice, but you can include more if necessary.
85     * You can include text in a `ipynb` file if you prefer.
86     * There is no penalty for submitting more or less text.
87     * There is no need to include references if you did not use existing resources. You don't have to reference the supporting material or the lab notes.

87 If you use external resources, e.g. code from online sources, then provide a reference in the report.
88     * You can run your scripts in any Python IDE you prefer, including Visual Studio Code, Colab, etc.
89
90     **[20 marks]**
91
92 **7.** Submit quality scripts following coding practices such as functional or object-oriented programming, separate Python files to break functionality for readability, use of `main`, and other standard practices. Feel free to use Jupyter Notebooks, but provide clear and well-organized code.
93
94
95     **[20 marks]**
96
97
98
99 ### Part 2: Supporting material
100
101 The following supporting scripts can help you kickstart your project. You spend some time familiarising yourself with the suggested libraries and scripts. You can use any library you prefer, not necessarily the ones provided below.
102
103 * Install the required libraries using the following command.
104
105 ```bash
106 pip3 install pytube moviepy speechrecognition spacy nltk NRCLex Path TextBlob
107 ```
108
109 > [!TIP]
110 >
111 > Depending on your Python installation, you might need to use `pip3` instead of `pip` .
112
113 * Explore the following script that downloads a

```
113 video in a `video_output folder.
114
115 ```python
116 from pytube import YouTube
117
118 url = "https://www.youtube.com/watch?v=jNQXAC9IVRw"
119 yt = YouTube(url)
120 stream = yt.streams.get_highest_resolution()
121 print(f"Downloading video: {yt.title}")
122 stream.download(output_path="video_output")
123 print(f"Download completed: {yt.title}")
124 ```
125
126 > The output is as follows:
127
128 ```bash
129 Downloading video: Me at the zoo
130 Download completed: Me at the zoo
131 ```
132
133 > [!NOTE]
134 >
135 > The files are stored as `mp4`; you can use the
    following script to create a folder.
136 >
137 > ```python
138 > from pathlib import Path
139 > Path("output").mkdir(parents=True, exist_ok=True)
140 > ```
141
142 * The following script extracts the audio file `wav
    ` from a video clip `mp4`.
143
144 ```python
145 import moviepy.editor as mp
146
147 video = mp.VideoFileClip("Me at the zoo.mp4")
148   video.audio.write_audiofile("Me at the zoo.wav")
149 ```
150
151 > [!NOTE]
```

```
152 >
153 > The script extracts the audio files in `wav`
    format in your local directory.
154
155 * If you like to store outputs in a folder you can
    use the following script:
156
157 ```python
158 import os
159
160 output_folder = "output"
161 output_filename = "myfile.txt"
162 output_folder = os.path.join(output_folder,
    output_filename)
163 ```
164
165 * The following script extracts the text from the
    audio file.
166
167 ```python
168 import speech_recognition as sr
169 recognizer = sr.Recognizer()
170 with sr.AudioFile("Me at the zoo.wav") as source:
171     audio = recognizer.record(source)
172 text = recognizer.recognize_google(audio)
173 print(text)
174 ```
175
176 > [!NOTE]
177 >
178 > The scripts extract the audio as a string. The
    following script saves the string to a text file.
179 >
180 > ```python
181 > file_path = 'Me at the zoo.txt'
182 > with open(file_path, 'w') as file:
183 >     file.write(text)
184 > print(f'Text has been written to {file_path}')
185 > ```
186
187 * The following script extracts the sentiment from a
```

```
187   text file.
188
189 ```python
190 from textblob import TextBlob
191
192 text = "all right so here we are one of the
      elephants and cool thing about these guys today is
      that they have really really really really really
      really long trunks and that's that's cool"
193
194 blob = TextBlob(text)
195 print(blob.sentiment)
196 print(blob.sentiment.polarity)
197 print(blob.sentiment.subjectivity)
198 ```
199
200 > [!NOTE]
201 >
202 > * Polarity measures the text's positive or
      negative tone. It ranges from `-1.0` (very negative
      ) to `1.0` (very positive). A value of `0` indicates
       a neutral sentiment.
203 >
204 > * Subjectivity measures how subjective or
      objective the text is. It ranges from `0.0` (very
      objective) to `1.0` (very subjective). A subjective
      text contains personal opinions, emotions, or
      judgments, whereas an objective text is based on
      factual information.
205
206 * Extract the video subtitles in Spanish - or any
      other language you prefer.
207
208 ```python
209 from textblob import TextBlob
210
211 text = "all right so here we are one of the
      elephants and cool thing about these guys today is
      that they have really really really really really
      really long trunks and that's that's cool"
212
```

```
213 blob = TextBlob(text)
214 blob_translated = blob.translate(from_lang='en', to
    ='es')
215 print(blob_translated)
216 ```
217
```

218 > [!NOTE]
219 >
220 > The output will be stored in the `blob_translated` in Spanish.

221

222 * Extract the emotions of the text using the `spacy`, `nltk` and the `NRCLex` libraries.

223

```python
225 import spacy,nltk
226 from nrclex import NRCLex
227 nlp = spacy.load('en_core_web_sm')
228 nltk.download('punkt')
229
230 text = "all right so here we are one of the
    elephants and cool thing about these guys today is
    that they have really really really really really
    really long trunks and that's that's cool"
231
232 doc = nlp(text)
233
234 full_text = ' '.join([sent.text for sent in doc.
    sents])
235
236 emotion = NRCLex(text)
237
238 print("Detected Emotions and Frequencies:")
239 print(emotion.affect_frequencies)
240 ```
241
```

242 > [!NOTE]
243 >
244 > The output will be the video emotions in a dictionary, such as:
245 >

```
246 > ```python
247 > Detected Emotions and Frequencies:
248 > {'fear': 0.0, 'anger': 0.0, 'anticip': 0.0, 'trust
    ': 0.0, 'surprise': 0.0, 'positive': 0.
    6666666666666666, 'negative': 0.0, 'sadness': 0.0, '
    disgust': 0.0, 'joy': 0.0, 'anticipation': 0.
    3333333333333333}
249 > ```
250
```