

```

1  import pandas as pd
2  import time
3  from joblib import Parallel, delayed
4  from datetime import datetime
5  import os
6  import shutil
7  import logging
8  from multiprocessing import cpu_count
9  from openpyxl import Workbook
10 from openpyxl.styles import PatternFill, Border, Side
11 from openpyxl.utils import get_column_letter
12 from concurrent.futures import ThreadPoolExecutor
13
14 # Set up logging
15 logging.basicConfig(filename='process_log.log', level=logging.INFO)
16
17 # Start the timer for the total process
18 start_time = time.time()
19
20 # Define the file paths to search
21 file_paths = [
22     r"C:\Users\Stephen\Documents\1. A Work\Test area\Folder 1\large_test_file_1.xlsx",
23     r"C:\Users\Stephen\Documents\1. A Work\Test area\Folder 1\large_test_file_2.xlsx",
24     r"C:\Users\Stephen\Documents\1. A Work\Test area\Folder 1\large_test_file_3.xlsx",
25     # Add more file paths as needed (up to 16 files)
26 ]
27
28 # Define the local directory where the files will be copied to
29 local_directory = r"C:\Users\Stephen\Documents\1. A
Work\python-full-3.12.5-windows-x86_64\python-full-3.12.5-windows-x86_64\Scripts"
30
31 # Function to copy a single file
32 def copy_file(file_path, local_directory):
33     file_name = os.path.basename(file_path)
34     local_file_path = os.path.join(local_directory, file_name)
35     shutil.copy2(file_path, local_file_path)
36     print(f"Copied {file_name} to {local_file_path}")
37     return local_file_path
38
39 # Function to copy files in parallel using ThreadPoolExecutor
40 def copy_files_to_local(shared_file_paths, local_directory, max_workers=8):
41     if not os.path.exists(local_directory):
42         os.makedirs(local_directory)
43
44     # Use ThreadPoolExecutor to copy files concurrently
45     local_file_paths = []
46     with ThreadPoolExecutor(max_workers=max_workers) as executor:
47         futures = [executor.submit(copy_file, file_path, local_directory) for file_path
in shared_file_paths]
48         for future in futures:
49             local_file_paths.append(future.result())
50
51     return local_file_paths
52
53 # Copy files using parallel threads
54 local_file_paths = copy_files_to_local(file_paths, local_directory, max_workers=8)
55
56 # Load the Excel file containing the search terms
57 search_terms_df = pd.read_excel(r"C:\Users\Stephen\Documents\Extracted_Info.xlsx")
58
59 # Extract the second row (index 0 in DataFrame)
60 row = search_terms_df.iloc[0]
61
62 # Function to normalize DOB to a consistent format (dd/mm/yyyy)
63 def normalize_dob(date_str):
64     try:

```

```

65         return datetime.strptime(date_str.strip(), "%d/%m/%Y").strftime("%d/%m/%Y")
66     except ValueError:
67         return date_str.strip()
68
69 # Function to clean strings
70 def clean_string(text):
71     return text.replace('\n', ' ').replace('\t', ' ').strip()
72
73 # Create a dictionary with the extracted search terms
74 new_search_terms = {
75     "first_name": clean_string(str(row["First Name"])),
76     "last_name": clean_string(str(row["Last Name"])),
77     "id": clean_string(str(row["ID 1"])),
78     "dob": normalize_dob(clean_string(str(row["DOB"]))),
79     "ac_no": clean_string(str(row["AC No."])),
80     "phone": clean_string(str(row["Phone"]))
81 }
82
83 # Define the search terms by file and sheet
84 search_terms_per_file = {
85     local_file_paths[0]: {
86         "Sheet1": new_search_terms,
87         "Sheet2": new_search_terms,
88         "Sheet3": new_search_terms,
89     },
90     local_file_paths[1]: {
91         "Sheet1": new_search_terms,
92     },
93     local_file_paths[2]: {
94         "Sheet1": new_search_terms,
95     },
96 }
97
98 # Function to search for terms in a specific Excel file
99 def process_excel_file(file_path, search_terms_per_file):
100     all_results = []
101     try:
102         excel_data = pd.read_excel(file_path, sheet_name=None)
103         for sheet_name, sheet_data in excel_data.items():
104             cleaned_search_terms = {term_name: clean_string(str(term_value).lower()) for
105                                     term_name, term_value in search_terms_per_file[file_path].get(sheet_name,
106                                                     {}).items()}
107             headers = list(sheet_data.columns)
108             for index, row in sheet_data.iterrows():
109                 found_terms = set()
110                 matched_values = {}
111                 for term_name, term_value_str in cleaned_search_terms.items():
112                     if any(term_value_str in clean_string(str(cell).lower()) for cell in
113                           row):
114                         found_terms.add(term_name)
115                         matched_values[term_name] = term_value_str
116                 if len(found_terms) >= 2:
117                     result = {
118                         "Matched Terms": matched_values,
119                         "File": file_path,
120                         "File Name": os.path.basename(file_path), # Adding file name
121                         "Sheet": sheet_name,
122                         "Row Number": index + 1,
123                         "Headers": headers,
124                         "Row Data": row.to_dict()
125                     }
126                     all_results.append(result)
127     except Exception as e:
128         logging.error(f"Error processing {file_path}: {e}")
129     return all_results

```

```

127
128 # Parallel processing to handle files in batches
129 def process_files_in_parallel(file_paths, search_terms_per_file):
130     all_results = Parallel(n_jobs=cpu_count() - 1)(delayed(process_excel_file)(file_path,
131         search_terms_per_file) for file_path in file_paths)
132     return [item for sublist in all_results for item in sublist]
133
134 # Execute the parallel processing
135 all_results = process_files_in_parallel(local_file_paths, search_terms_per_file)
136
137 # Prepare for writing to Excel
138 workbook = Workbook()
139 sheet = workbook.active
140 sheet.title = "Results"
141
142 # Add headers for file, sheet, and row
143 sheet.append(["File", "Sheet", "Row", "-----", "Row Data"])
144
145 # Define styles
146 header_fill = PatternFill(start_color="0033A0", end_color="0033A0", fill_type="solid")
147 # Dark blue for headers
148 data_fill = PatternFill(start_color="D3D3D3", end_color="D3D3D3", fill_type="solid")
149 # Light gray for data
150
151 # Define borders
152 thin_border = Border(left=Side(style='thin'), right=Side(style='thin'), top=Side(style=
153 'thin'), bottom=Side(style='thin'))
154
155 # Freeze top row for better navigation
156 sheet.freeze_panes = "A2"
157
158 # Write data with formatting, starting from Column D, without "Headers" or "Row Data"
159 labels
160 for idx, result in enumerate(all_results):
161     row_fill = header_fill if idx % 2 == 0 else data_fill
162
163     # Write the actual headers and data without labels
164     row_num = sheet.max_row + 1
165     sheet.append([result["File Name"], result["Sheet"], result["Row Number"], "" ] +
166         result["Headers"])
167
168     # Add hyperlink to the file path
169     sheet.cell(row=row_num, column=1).hyperlink = result["File"] # Adds hyperlink to
170     cell with the file name
171     sheet.cell(row=row_num, column=1).value = result["File Name"] # Ensures cell shows
172     the file name
173
174     for col_num in range(5, len(result["Headers"]) + 5):
175         cell = sheet.cell(row=row_num, column=col_num)
176         cell.fill = header_fill
177         cell.border = thin_border # Add borders
178
179     # Write the data row starting from Column D
180     row_num = sheet.max_row + 1
181     sheet.append(["", "", "", ""] + list(result["Row Data"].values()))
182     for col_num in range(5, len(result["Row Data"].values()) + 5):
183         cell = sheet.cell(row=row_num, column=col_num)
184         cell.fill = data_fill
185         cell.border = thin_border # Add borders
186
187 # Adjust column widths for better readability
188 for col in sheet.columns:
189     max_length = 0
190     column = col[0].column_letter
191     for cell in col:
192         try:

```

```

185         if len(str(cell.value)) > max_length:
186             max_length = len(cell.value)
187     except:
188         pass
189     adjusted_width = (max_length + 2) if max_length < 30 else 30
190     sheet.column_dimensions[column].width = adjusted_width
191
192 # Save the workbook
193 current_time = datetime.now().strftime("%Y-%m-%d %H-%M-%S")
194 output_file = f"{current_time} Search_results.xlsx"
195 workbook.save(output_file)
196
197 # Generate text file with output file name
198 def generate_text_file(output_file_name):
199     text_file_path = r"C:\Users\Stephen\Documents\1. A
200     Work\python-full-3.12.5-windows-x86_64\python-full-3.12.5-windows-x86_64\Scripts\gene
201     rated_file_name.txt"
202     with open(text_file_path, 'w') as text_file:
203         text_file.write(f"{output_file_name}\n")
204     print(f"Text file generated at: {text_file_path}")
205
206 generate_text_file(output_file)
207
208 # End the timer and calculate the total runtime
209 end_time = time.time()
210 total_runtime = end_time - start_time
211
212 print(f"Total search complete. Results saved to {output_file}")
213 print(f"Total Runtime: {total_runtime:.2f} seconds")

```