

# Coursework Coversheet

School of Geography. FACULTY OF ENVIRONMENT  
GEOG5403 Creative Coding for Urban Problems



## UNIVERSITY OF LEEDS

Student ID	201788236/201787698/201790263/201 776257/201793166/201788996	Mark  Less deduction (state reason)  Final Mark	
Group Number	Group 5		
Assignment title	Child labour in global supply chain		
Marker			

Justification of mark (using specific text from criteria):

To improve your work for next time:

1.

2.

3.

Additional comments:

Group 5

# Child labour in global supply chain

Alexander Thornton, Shenghang Yuan, Subaru Shimizu, Yuxi Liu, Zhaojing Wang, Zhihao Zhang



## Background

Nearly 1 in 10 children are subjected to child labour worldwide,  
With some forced into hazardous work through trafficking.

### Why do families choose children to be employed

- Economic hardship
- Migrant and refugee children

### Why companies use child labour

- High Demand, Low Skill
- Lower price
- Less Oversight



## Background

Nearly 1 in 10 children are subjected to child labour worldwide,  
With some forced into hazardous work through trafficking.

### Why child labour data is not available/difficult to obtain

- Informal sectors
- Hidden or illegal contexts
- Reluctance to report

### Consequence

- Extreme bodily and mental harm
- Compound social inequality and discrimination





# Approach to the challenge and Used Data

1. **Challenge:** It is difficult to understand the current situation of child labour due to the data availability as companies tend to hide negative information.
2. **Objective:** To understand where, in which commodity production and by which company the child labour is executed and explore ways to solve these problems.
3. **Approach:** Use the news article provided by HACE, and then extract articles that contain our keywords, extract information using ChatGPT API, and finally visualize the results.
4. **Used dataset:** Text files of news articles by HACE, we mainly selected dataset from different countries, companies and commodities to analyze and see the relationship between them.



# Flow of Analysis

- 1. Identification of keywords:** Keywords related to child labour and commodities were identified. These keywords were used in following steps to filter and extract relevant news articles.
- 2. Reading articles:** These articles were obtained from some data source, for example a news website.
- 3. Extracting articles that contain our keywords:** We identified each text file to figure out whether it contains keywords (child labour). If so, this file was extracted for analysis.
- 4. Extracting key information using ChatGPT API:** We used ChatGPT to process and analyse text files to obtain information, including countries, commodities and companies involved.
- 5. Visualisation (including minor data cleaning):** The extracted information was used for visualisation. It includes the use of interactive maps to show the distribution of countries, and the display of information about commodities and companies.

# Methods

## Challenge 1: Massive dataset that contains news articles irrelevant to child labour

- The provided dataset has 4,561 news article files, which is too massive, and contains ones irrelevant to child labour.

## Solution 1: Extracted news articles that contain child labour-related keywords

- Prepared a list of 61 child-labour-related keywords with the assistance of ChatGPT.
  - Examples: *Child labour, Child labor, Kid worker, Child exploitation, Child trafficking, Child slavery, Child bondage, Juvenile exploitation, etc.*
- Extracted news articles that contain child labour-related keywords.

Please provide a list of keywords in relation to child labour.



# Methods

## Challenge 2: Difficulty in extracting key information from news articles.

Example:



child labortaints production of batteries for electric carmakers, amnesty says news today's news us politics world tech science weather opinion originals the 360 skullduggery podcast conspiracyland finance my portfolio yahoofinance plus news screeners markets videos watchlists personalfinance crypto industries sports fantasy dailyfantasy nfl nba mlb nhl ncaaf ncaab more ncaaw mma sportsbook soccer tennis nascar golf boxing cycling wnba usfl indycar horse racing olympics gamechannel rivals podcasts videos rss team apparel and gear shop breakingt shirts newsfeed entertainment celebrity tv movies music life health style and beauty parenting horoscopes wel-being food travel terms privacy dashboard feedback © 2023 all rights reserved. about our ads advertising careers

yahoo news search query news finance sports more news today's news us politics world covid-19 climate change health science originals life health parenting style and beauty horoscopes shopping entertainment role recall anniversary party game changers the it list under the covers the neverweres are the kids alright? how to watch interviews videos shopping finance my portfolio watchlists markets news videos yahoo finance plus screeners personal finance crypto industries sportsfantasy nfl nba mlb nhl college football college basketball soccer mma yahoo sports am editions us englishus españaaustralia englishcanada englishcanada françaisfrance françaisgermany deutschhong kong 繁體中文malaysia englishsingapore englishtaiwan 繁體中文uk english mail sign in advertisementclose this contentreuterschild labortaints production of batteries for electric carmakers, amnesty saysby lin taylor september 30, 2016 at 3:38 amlink copiedread full articleby lin taylor london (thomson reuters foundation) - leading electric carmakers may be unwittingly using child labor to produce batteries for vehicles that have grownin popularity for using clean energy aimed at limiting global warming, amnesty international said onfriday. the human rights watchdog said cobalt used in lithium ion batteries for electric vehicles, phones and laptops could come from mines in democratic republic of congo (drc) that use child labor. it accused carmakers including gm, renault-nissan, fiat chrysler, volkswagen, daimler and tesla offailing to map the supply of cobalt from mines in congo to smelters and on to battery-makers. as a result, electric cars sold across the globe could contain traces of the metal produced each year by informal congolese mines without companies knowing, the group said. "we have found that there is significant risk that cobalt mined by children could be entering their supply chains," said markdummett, human rights researcher at amnesty. due to the growing demand for lithium ion batteries in electric cars and consumer goods like mobile phones, dummett said carmakers had an obligation to trace their supply chain to ensure children were not used to mine cobalt. "frankly companies owe it to their consumer to be transparent about their supplies and to map out their supply chains so that they know where its coming from," he said in an interview with the thomson reuters foundation. more than half of the world's cobalt comes fromdrc, amnesty said, and 20 percent of the mineral was mined by hand. millions of congolese work in informal mining, with rudimentary tools and usually without legal

This article mentions the possibility that cobalt mined by children in DRC is supplied to EV manufacturers.

### Key information

- **Country:** Democratic Republic of Congo (DRC)
- **Company:** GM, Renault-Nissan, Fiat Chrysler, Volkswagen, Daimler, Tesla
- **Commodity:** Cobalt



Source of figure: [CNN](#)

# Methods

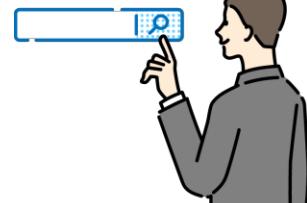
## Challenge 2: Difficulty in extracting key information from news articles.

### Solution 2: Extracted key information using ChatGPT API.

- Utilised ChatGPT API to extract key information from news articles.
- Provided ChatGPT with a list of 62 relevant commodities to improve returned answers.

Can you find the **companies** that are mentioned in relation to child labour exploitation?

Here is an example of your message: 'Microsoft, Alphabet'.



Can you find the **countries** that are mentioned in relation to child labour exploitation?

Here is an example of your message: 'United States, Indonesia'.

Can you find the **commodities** that are mentioned in relation to child labour exploitation?

Please provide a list of commodity names, using the names in the commodity list:

'{commodity\_labour}'.

Here is an example of your message: 'Jewelry, Apparel'.

Relevant commodities: Gold, Cobalt, Textiles, Apparel, Footwear, Jewellery, Batteries, etc.



# Result Description

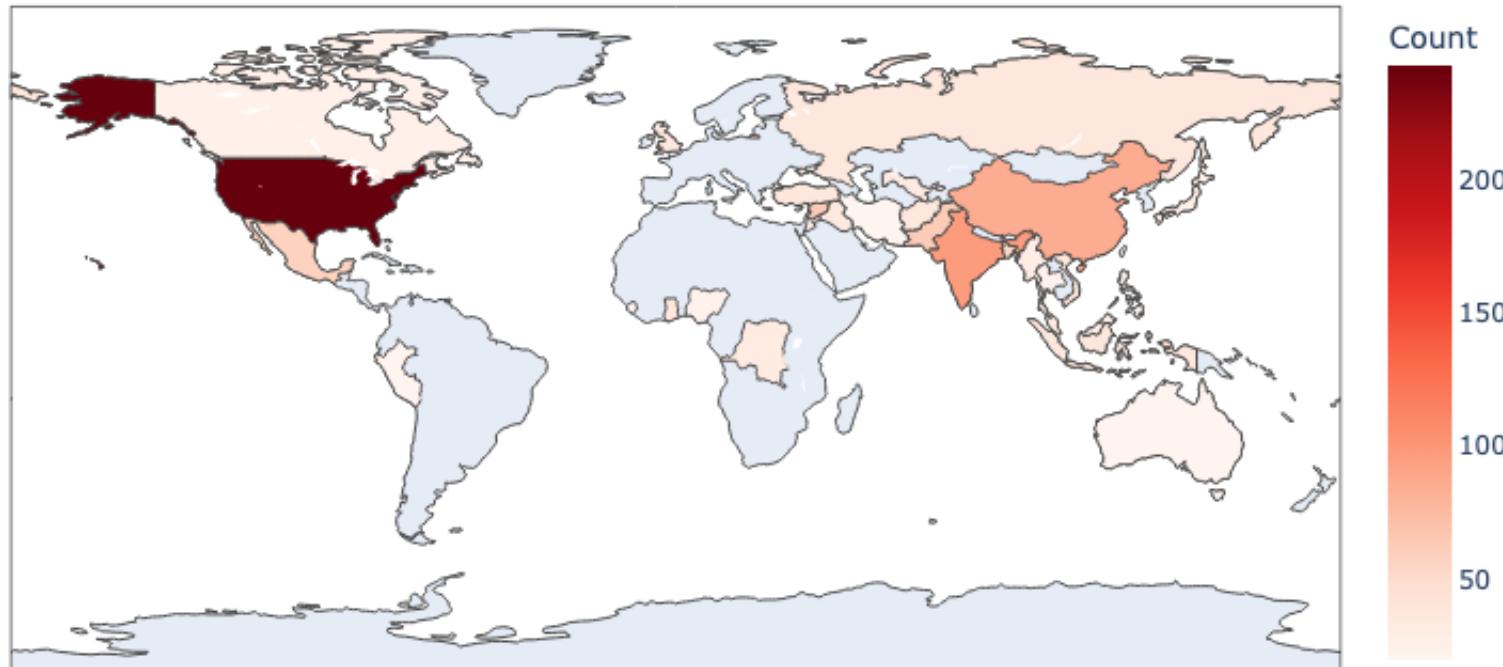
Summarised returned answers from ChatGPT.

- Company
- Commodity
- Country

Company	Commodity	Country
None.	None	Brazil
Apple, Tesla	Cobalt	Democratic Republic of Congo
Apple Inc, Sony Corp, Volvo	Cobalt	Democratic Republic of Congo
None	Cotton	None
None	'Cobalt'	Canada, Democratic Republic of Congo
None	Child Labor, Clothing, Footwear, Toys	United States, Germany, Japan, Nicaragua, Mexico, United Arab Emirates, Kuwait, Oman, Georgia, Florida, Alabama, Alaska, Tennessee, Indiana, Oklahoma, Louisiana
None	Gravel, Construction materials, Stone, Fish, Fish processing	Myanmar
None	None	United States
None	Cotton	Turkmenistan
None	Toys, Handicrafts, Textiles, Apparel, Carpentry, Garments, Pharmaceuticals.	Bolivia
None	None	Brazil, United States
None	Fish, Coffee, Sugarcane, Cotton, Tobacco, Alcohol, Leather, Cattle	China, India, Bangladesh, Brazil, Indonesia, Philippines, Vietnam, Thailand, Egypt, Cambodia, Pakistan, Malaysia, Myanmar, Turkey, Mexico, Nigeria, Ghana, Kenya, El Salvador, Guatemala, Honduras, Nicaragua, Panama, Colombia, Costa Rica, Dominican Republic, Haiti, Jamaica
None	None	None.
cez, korea electric power (kepcos), norges bank investment management (nbim)	Coal, Textiles	Turkey
None	Agricultural produce, Textiles, Apparel, Construction materials	Vietnam
None	None	Jordan, United States, Syria, Iraq, Yemen, Libya, Pakistan
None	None	Chile, Yemen, Syria, Iraq, Jordan, the Netherlands, Monaco, East Timor, Australia, India
None	Toys	China

## Visualisations

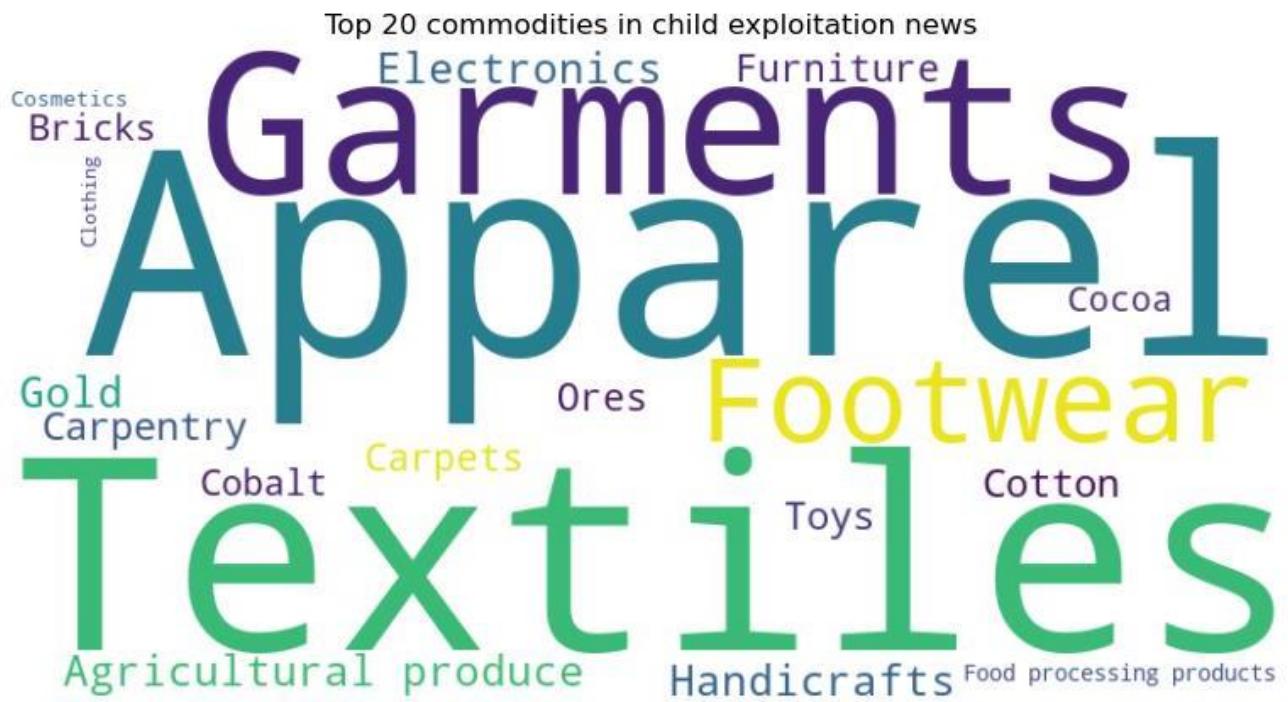
Frequencies of child exploitation news by country (over 20 counts only)



Country	Count
United States	243
India	97
China	86
Syria	64
Pakistan	63
Mexico	58
Bangladesh	54
Afghanistan	39
Ghana	39
Russia	39
Indonesia	38
United Kingdom	38
Vietnam	38
Turkey	37
Iraq	36
Uzbekistan	35
Democratic Republic of Congo	33
Myanmar	30
Jordan	30
Malaysia	29

## Visualisations

### Commodities:

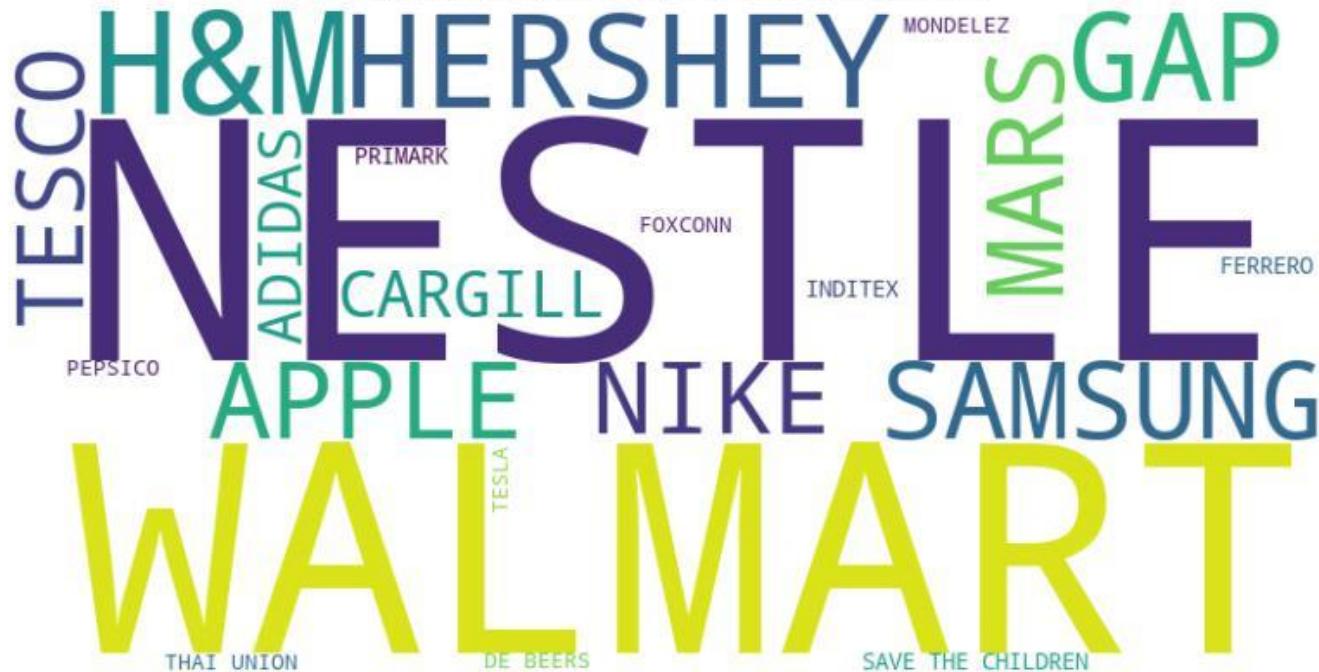


Commodity	Count
Apparel	115
Textiles	101
Garments	61
Footwear	57
Agricultural produce	45
Electronics	44
Handicrafts	43
Gold	43
Cotton	37
Toys	37
Carpentry	36
Ores	32
Carpets	31
Furniture	30
Cobalt	28
Bricks	27
Cocoa	25
Food processing products	25
Cosmetics	22
Clothing	18

## Visualisations

### Companies:

Top 20 companies in child exploitation news



Company	Count
NESTLE	20
WALMART	15
H&M	14
HERSCHEY	11
GAP	10
SAMSUNG	9
APPLE	9
MARS	9
NIKE	8
TESCO	7
CARGILL	7
ADIDAS	6
SAVE THE CHILDREN	6
PRIMARK	6
FERRERO	6
INDITEX	6
FOXCONN	5
PEPSICO	5
THAI UNION	5
TESLA	5
MONDELEZ	5
DE BEERS	5

## Findings from Analysis

### Case 1: United States

Company and Commodity

- **Driscoll's and Walmart:** Agricultural products such as berries and vegetables
- **Coca-Cola and other beverage manufacturers:** the production of sugar and other agricultural products, especially cocoa, fruits and vegetables.

Insights from articles

- In agriculture and textile manufacturing, companies may exploit child labour to perform **low-skilled and labour-intensive work** in order to reduce costs. In agriculture, in particular, child labour is easier to overlook because production sites are scattered and difficult to regulate.
- Farms in the United States rely on **the immigrant labour force (Mexican workers)** and exploit these people.

### Case 2: China

Company and Commodity

- **Electronics and clothing/textile** are two major commodities associated with child labour, involving several internationally renowned brands such as **Foxconn, Apple, Gap and H&M**.

## Solutions and Policies

### Strengthening regulation and enforcement

<https://www.google.com>



- Strengthen regulations to ensure that enterprises, including their suppliers and subcontractors in developing countries, comply with international standard labour requirements.
- Strengthen enforcement to punish businesses that violate labour laws.
- Encourage companies to inspect their supply chains to ensure that child labour is not involved.

### Improve education and economic conditions

- Provide more educational resources such as free textbooks and uniforms to ensure more children have access to good education.
- Provide financial support and job opportunities to those who have children.

### Create social co-operation

- Work together to address child labour through cooperation between government, business and the community. Establish a social support system to provide necessary assistance to children exploited by child labour.



# Future work

## Preparing readable dataset

- Cleaning/reformatting the dataset by removing unnecessary texts such as headers/footers will improve readability.

## Improving prompts given to ChatGPT for better accuracy

- Providing ChatGPT with a distinct prompt, including background knowledge and specific areas of concern, to return our expected responses.

## Developing natural language processing models for sustainability and accountability

- Since relying on ChatGPT API can be costly, advanced natural language processing models (e.g., entity recognition systems) that improve accuracy and accountability will be recommended.

## Improving output data cleaning and validation process

- Responses from ChatGPT require manual data cleaning, such as unifying names (e.g. Apple, Apple Inc.). Giving a list of companies or commodities will improve the output data.
- Due to the complex texts, the output needs validation to ensure accuracy.



## References

- Anon n.d. *Child labor versus educational attainment Some evidence from Latin America Original*.
- Edmonds, E.V. and Schady, N. 2012. Poverty alleviation and child labor. *American Economic Journal. Economic Policy*. 4(4), pp.100–124.
- Humphries, J. 2013. Childhood and child labour in the British industrial revolution1: CHILDHOOD AND CHILD LABOUR. *The economic history review*. 66(2), pp.395–418.
- Nieuwenhuys, O. 1996. The paradox of child labor and anthropology. *Annual review of anthropology*. 25(1), pp.237–251.
- Ray, R. 2000. Child labor, child schooling, and their interaction with adult labor: Empirical evidence for Peru and Pakistan. *The World Bank economic review*. 14(2), pp.347–367.
- Zapata, D., Contreras, D. and Kruger, D. 2011. Child labor and schooling in Bolivia: Who's falling behind? The roles of domestic work, gender, and ethnicity. *World development*. 39(4), pp.588–599.



# Contribution

**Coding:** All group numbers

**Background and Approach:** Shenghang Yuan

**Data description and Methods:** Alexander Thornton, Zhihao Zhang

**Visualisation and Future work:** Subaru Shimizu

**Slides creation and design:** Zhaojing Wang, Yuxi Liu

**Speech:** Alexander Thornton, Subaru Shimizu, Zhaojing Wang



# Thanks

```
In [2]: 1 import zipfile
2 import os
3
4 def extract_zip(zip_file, extract_to):
5     with zipfile.ZipFile(zip_file, 'r') as zip_ref:
6         zip_ref.extractall(extract_to)
7
8 ## Example usage
9 # zip_file = 'OneDrive_2024-03-19.zip' # Replace 'example.zip' with the name of your zip file
10 # extract_to = r'C:\Users\23510\Creative Coding for Urban Problems\Week 4' # Replace 'extracted_files' with the directory where you want to extract the files
11
12 ## Create the directory if it doesn't exist
13 # os.makedirs(extract_to, exist_ok=True)
14
15 ## Extract the zip file
16 # extract_zip(zip_file, extract_to)
17
```

```
In [24]: 1 import os
2
3 import pandas as pd
4 import numpy as np
5 import seaborn as sns
6 palette = 'bwr'
7 sns.set_palette(palette)
8
9 import matplotlib.pyplot as plt
10 from matplotlib.pyplot import figure
11 pd.options.mode.chained_assignment = None
12
13 #import datetime
14 from datetime import datetime
15
16 import plotly.express as px
17 import plotly.graph_objects as go
18 import plotly.io as pio
19 import plotly.offline as pyo
20
21 %matplotlib inline
22
23 from pysal.lib import weights
24 from pysal.lib import io
25
26 import contextily as cx
27 import geopandas as gpd
28 from shapely.geometry import Point
29 from shapely import wkt
30
31 import folium
32 from folium import plugins
33
34 import spacy
35 import openai
36 import requests
37
38 from wordcloud import WordCloud
39 import matplotlib.pyplot as plt
40
41 pd.set_option('display.max_columns', 200)
42 pd.set_option('display.max_rows', 200)
```

```
In [14]: 1 # Assign directory
2
3 directory = r'/Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02_Sem2/5403_Creative Coding/Week4/Alex/dataset'
4
5 # get all files under directory
6 file_list = os.listdir(directory)
7
8 # change the extention to .txt
9 for file_name in file_list:
10     # file's full path
11     file_path = os.path.join(directory, file_name)
12     # change the extention
13     new_file_path = os.path.splitext(file_path)[0] + '.txt'
14     try:
15         # change the file name
16         os.rename(file_path, new_file_path)
17         print("Renamed [file_path] to [new_file_path]")
18     except Exception as e:
19         print("Error renaming file {file_path}: {e}")
```

e/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20170123191500-894.txt  
Renamed /Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20160604200000-T2680.txt to /Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20160604200000-T2680.txt  
Renamed /Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20161008091500-484.txt to /Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20161008091500-484.txt  
Renamed /Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20161118131500-1982.txt to /Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20161118131500-1982.txt  
Renamed /Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20151201211500-2747.txt to /Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20151201211500-2747.txt  
Renamed /Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20161125101500-225.txt to /Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20161125101500-225.txt  
Renamed /Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20170615204500-612.txt to /Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20170615204500-612.txt  
Renamed /Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20161124071500-T703.txt to /Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20161124071500-T703.txt  
Renamed /Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20160609040000-247.txt to /Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20160609040000-247.txt  
Renamed /Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20160612194500-T922.txt to /Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20160612194500-T922.txt  
Renamed /Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20170503060000-T556.txt to /Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20170503060000-T556.txt

## Step 1: Preparation

### 1-1: Read Keywords of Child Labour and Commodities

```
In [15]: 1 # Child labour related keywords
2 keywords_excel = pd.read_excel('key_words.xlsx')
3 keywords_excel.head()
```

	Relevant to Child labour	Relevant to commodities	Fortune 500
0	Child labour	Gold	Walmart Inc.
1	child labor	Cobolt	Amazon.com Inc.
2	kid worker	Textiles	Apple Inc.
3	Child Exploitation	Apparel	Exxon Mobil Corporation
4	Underage Labour	Footwear	Berkshire Hathaway Inc.

In [16]:

```
1 keyword_labour = []
2 for words in keywords_excel['Relevant to Child labour']:
3     keyword_labour.append(words)
4
5 keyword_labour = [keyword for keyword in keyword_labour if pd.notnull(keyword)]
6 keywords = keyword_labour
7 keywords
```

Out[16]:

```
['Child labour',
 'child labor',
 'kid worker',
 'Child Exploitation',
 'Underage Labour',
 'Juvenile Labour',
 'Child Work',
 'Child Employment',
 'Minor Labour',
 'Youth Labour',
 'Adolescent Labour',
 'Child Servitude',
 'Child Trafficking',
 'Child Servitude',
 'Underage Employment',
 'Child Exploitation',
 'Minor Workforce',
 'Youth Exploitation',
 'Child Trafficking',
 'Juvenile Exploitation',
 'Child Bondage',
 'Child Forced Labour',
 'Child Slavery',
 'Child Indentured Labour',
 'Child Sweatshop Labour',
 'Child Exploitative Labour',
 'Child Labor Trafficking',
 'Child Commercial Exploitation',
 'Child Abuse of Labour',
 'Child Employment Exploitation',
 'Child Labor Exploitation',
 'Child Labor Abuse',
 'Child Labor Trafficking',
 'Child Work Exploitation',
 'Child Work Abuse',
 'Child Labor Violation',
 'Child Workforce Exploitation',
 'Child Labor Malpractice',
 'Child Labor Infringement',
 'Child Workforce Injustice',
 'Child Labor Iniquity',
 'Child Labor Misuse',
 'Child Labor Maltreatment',
 'Child Abuse of Labour',
 'Child Sexual Exploitation',
 'Child Labor Exploitation',
 'Child Labor Abuse',
 'Child Pornography',
 'Child Labor Trafficking',
 'Child Labor Violation',
 'Child Labor Infringement',
 'Child Prostitution',
 'Child Marriage',
 'Child Sex Trafficking',
 'Child Labor Misuse',
 'Child Labor Maltreatment',
 'Child Labor Iniquity',
 'Child Commercial Exploitation',
 'Child Soldiering',
 'Child Domestic Servitude',
 'Child Labor in Agriculture']
```

In [17]:

```
1 # Commodity related keywords
2 commodity_labour = []
3 for words in keywords_excel['Relevant to commodities']:
4     commodity_labour.append(words)
5
6 commodity_labour = [keyword for keyword in commodity_labour if pd.notnull(keyword)]
7 commodity_labour
8
```

Out[17]:

```
['Gold',
'Cobolt\xao',
'Textiles',
'Apparel',
'Footwear',
'Toys',
'Handicrafts',
'Electronics',
'Cosmetics',
'Furniture',
'Food processing products',
'Ores',
'Bricks',
'Recycled materials',
'Paper products',
'Utensils',
'Medical supplies',
'Packaging materials',
'Agricultural produce',
'Metal products',
'Building materials',
'Leather goods',
'Carpets',
'Rugs',
'Carpentry',
'Garments',
'Watches',
'Jewelry',
'Glassware',
'Ceramics',
'Pottery',
'Tiles',
'Plastics',
'Paintings',
'Sculptures',
'Musical instruments',
'Watches',
'Eyewear',
'Cosmetics',
'Pharmaceuticals',
'Automotive parts',
'Sporting goods',
'Fireworks',
'Matches',
'Cigarettes',
'Handbags',
'Wallets',
'Stationery',
'Carpentry tools',
'Construction materials',
'Cleaning products',
'Agricultural chemicals',
'Textile dyes',
'Glues and adhesives',
'Paints and coatings',
'Batteries',
'Electrical components',
'Plastics and rubber products',
'Carpets and rugs',
'Furniture upholstery',
'Disposable cutlery and dishes',
'Fishing gear']
```

## 1-2: Read Text (News Article) Files

In [18]:

```
1 def extract_text_file(file_path, encoding='utf-8'):
2     with open(file_path, 'r', encoding=encoding) as file:
3         text = file.read()
4     return text
5
6 # Example usage
7 file_path = r'/Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02_Sem2/5403_Creative Coding/Week4/Alex/dataset/20150310204500-2046.txt'
8 # Replace with the path to your text file
9 text_content = extract_text_file(file_path)
10 print(text_content)
11
```

new tax shocker for middle class to cater for the elderly, orphans - the standard

×

the standard group plc is a multi-media organization with investments in media platforms spanning newspaper print operations, television, radio broadcasting, digital and online services. the standard group is recognized as a leading multi-media house in kenya with a key influence in matters of national and international interest.

standard group plc hq office,  
the standard group center,mombasa road.  
p.o box 30080-00100,nairobi,kenya.  
telephone number: 0203222111, 0719012111  
email: [email protected]

news & current affairs

digital news

### 1-3: Find Text files related to Child labour

```
In [19]: 1  ## Find text files that contains the keywords.
2  def find_child_labor_references_in_folder(folder_path, keywords, max_files=10000, encoding='utf-8'):
3      child_labor_files = []
4
5      # Iterate over files in the folder
6      file_count = 0
7      for filename in os.listdir(folder_path):
8          if file_count >= max_files:
9              break
10
11      if filename.endswith('.txt'):
12          file_path = os.path.join(folder_path, filename)
13          with open(file_path, 'r', encoding=encoding, errors='ignore') as file:
14              text = file.read().lower() # Convert text to lowercase for case-insensitive matching
15
16      # Check for keywords in the text
17      for keyword in keywords:
18          if keyword in text:
19              child_labor_files.append(filename)
20              break # Once any keyword is found, no need to check further for this file
21
22      file_count += 1
23
24  return child_labor_files
25
26 # Example usage
27 folder_path = r'Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02_Sem2/5403_Creative Coding/Week4/Alex/dataset' # Replace with the path to your folder
28 child_labor_files = find_child_labor_references_in_folder(folder_path, keywords)
29
30 if child_labor_files:
31     print("Files mentioning child labor keywords:", child_labor_files)
32 else:
33     print("No files mentioning child labor keywords found.")
34
```

30074500-1584.txt', '20160825110000-2287.txt', '20180305224500-1251.txt', '20170612234500-1498.txt', '20170320201500-338.txt', '20170405181500-522.txt', '20160119223000-781.txt', '20151106171500-1092.txt', '20170211000000-1450.txt', '20180329100000-429.txt', '2016103214500-2066.txt', '20150713191500-1901.txt', '20170316043000-167.txt', '2018041233000-2775.txt', '20180425054500-1739.txt', '20160420161500-301.txt', '20161021190000-633.txt', '20160419201500-448.txt', '20151022224500-1531.txt', '20160614101500-1857.txt', '20160704191500-2891.txt', '20160706134500-2513.txt', '20151008004500-2494.txt', '20150905001500-92.txt', '20150829010000-542.txt', '20160621171500-1903.txt', '20180407230000-495.txt', '20180415034500-82.txt', '201703202111500-2508.txt', '20160114024500-168.txt', '20170428010000-2402.txt', '20170616051500-1786.txt', '20180418161500-652.txt', '20161129074500-718.txt', '20180313174500-278.txt', '20160315140000-2817.txt', '20160527061500-1056.txt', '20170802233000-2017.txt', '20170817060000-1186.txt', '20151205174500-1945.txt', '201603050121163000-1592.txt', '20170402200000-685.txt', '20170207084500-994.txt', '20160212160000-2058.txt', '20180302184500-2700.txt', '20170320071500-2049.txt', '20170331131500-1255.txt', '20170623014500-93.txt', '20151113161500-674.txt', '20170728133000-1361.txt', '20150424151500-1097.txt', '20161102043000-2335.txt', '2017042204500-636.txt', '2016020223000-550.txt', '20160710110000-1155.txt', '20180422051500-794.txt', '20170330180000-666.txt', '20160115080000-148.txt', '20150703054500-819.txt', '20170315120000-1536.txt', '20180314154500-2922.txt', '20160812054500-2455.txt', '20161224174500-507.txt', '20160920061500-2425.txt', '20170621163000-3.txt', '20170427111500-2058.txt', '20160205230000-284.txt', '20180418133000-1866.txt', '20170120173000-1073.txt', '20180328200000-2221.txt', '20170316160000-1788.txt', '20150501161500-2598.txt', '20170214033000-2022.txt', '20170315210000-177.txt', '20170322180000-1462.txt', '20170620030000-179.txt', '20160105114500-1354.txt', '20160422220000-809.txt', '20160613163000-1090.txt', '20170214170000-566.txt', '20170522223000-2506.txt', '20170722204500-1934.txt', '20170421134500-361.txt', '20180413134500-1401.txt', '20160503123000-427.txt', '20180403190000-1525.txt', '20170517211500-2346.txt', '20170614174500-804.txt', '20180404171500-654.txt', '20160808153000-2848.txt', '20160906164500-901.txt', '2016012223000-2899.txt', '2016042817300-0-1350.txt', '20180309040000-234.txt', '20170615203000-286.txt', '20161213070000-1956.txt', '20170323000000-2337.txt', '20170311063000-1680.txt', '20160725103000-2219.txt', '2017052100000-274.txt', '20150501003000-1733.txt', '20180416113000-2586.txt', '20160429044500-316.txt', '20180419201500-2398.txt', '20151218170000-1875.txt', '20170306111500-2058.txt', '201804031900-986.txt', '20170413204500-2870.txt', '20160630131500-558.txt', '20160409034500-1796.txt', '20170201014500-342.txt', '20170619153000-151.txt', '20161019054500-265.txt', '20170724164500-264.txt', '2017072020171500-811.txt', '20151011121500-510.txt', '20170308213000-281.txt', '20170505211500-1002.txt', '20170201173000-1922.txt', '20170624134500-1547.txt', '20180422164500-255.txt', '20170408053000-1355.txt', '20180417061500-1506.txt', '20170622193000-1329.txt', '20161017031500-1446.txt', '20180302204500-1892.txt', '2015110513000-1374.txt', '20160314083000-1134.txt', '2018030503000-789.txt', '20170501224500-33.txt', '20180329053000-806.txt', '20160927013000-83.txt', '20170612234500-885.txt', '20170316053000-1200.txt', '20170418133000-1200.txt', '201705211500-1200.txt', '20170612160000-1200.txt', '20170724164500-1200.txt', '20170825110000-1200.txt', '201709151500-1200.txt', '20171014160000-1200.txt', '20171113170000-1200.txt', '20171212180000-1200.txt', '20180111180000-1200.txt', '20180210190000-1200.txt', '20180319160000-1200.txt', '20180418170000-1200.txt', '20180517150000-1200.txt', '20180616140000-1200.txt', '20180715130000-1200.txt', '20180814120000-1200.txt', '20180913110000-1200.txt', '20181012100000-1200.txt', '20181111090000-1200.txt', '20181210180000-1200.txt', '20190109170000-1200.txt', '20190208160000-1200.txt', '20190307150000-1200.txt', '20190406140000-1200.txt', '20190505130000-1200.txt', '20190604120000-1200.txt', '20190703110000-1200.txt', '20190802100000-1200.txt', '20190901090000-1200.txt', '20191009080000-1200.txt', '20191108070000-1200.txt', '20191207060000-1200.txt', '20200106050000-1200.txt', '20200205040000-1200.txt', '20200304030000-1200.txt', '20200403020000-1200.txt', '20200502010000-1200.txt', '20200601000000-1200.txt', '20200701000000-1200.txt', '20200801000000-1200.txt', '20200901000000-1200.txt', '20201001000000-1200.txt', '20201101000000-1200.txt', '20201201000000-1200.txt', '20201301000000-1200.txt', '20201401000000-1200.txt', '20201501000000-1200.txt', '20201601000000-1200.txt', '20201701000000-1200.txt', '20201801000000-1200.txt', '20201901000000-1200.txt', '20202001000000-1200.txt', '20202101000000-1200.txt', '20202201000000-1200.txt', '20202301000000-1200.txt', '20202401000000-1200.txt', '20202501000000-1200.txt', '20202601000000-1200.txt', '20202701000000-1200.txt', '20202801000000-1200.txt', '20202901000000-1200.txt', '20203001000000-1200.txt', '20203101000000-1200.txt', '20203201000000-1200.txt', '20203301000000-1200.txt', '20203401000000-1200.txt', '20203501000000-1200.txt', '20203601000000-1200.txt', '20203701000000-1200.txt', '20203801000000-1200.txt', '20203901000000-1200.txt', '20204001000000-1200.txt', '20204101000000-1200.txt', '20204201000000-1200.txt', '20204301000000-1200.txt', '20204401000000-1200.txt', '20204501000000-1200.txt', '20204601000000-1200.txt', '20204701000000-1200.txt', '20204801000000-1200.txt', '20204901000000-1200.txt', '20205001000000-1200.txt', '20205101000000-1200.txt', '20205201000000-1200.txt', '20205301000000-1200.txt', '20205401000000-1200.txt', '20205501000000-1200.txt', '20205601000000-1200.txt', '20205701000000-1200.txt', '20205801000000-1200.txt', '20205901000000-1200.txt', '20206001000000-1200.txt', '20206101000000-1200.txt', '20206201000000-1200.txt', '20206301000000-1200.txt', '20206401000000-1200.txt', '20206501000000-1200.txt', '20206601000000-1200.txt', '20206701000000-1200.txt', '20206801000000-1200.txt', '20206901000000-1200.txt', '20207001000000-1200.txt', '20207101000000-1200.txt', '20207201000000-1200.txt', '20207301000000-1200.txt', '20207401000000-1200.txt', '20207501000000-1200.txt', '20207601000000-1200.txt', '20207701000000-1200.txt', '20207801000000-1200.txt', '20207901000000-1200.txt', '20208001000000-1200.txt', '20208101000000-1200.txt', '20208201000000-1200.txt', '20208301000000-1200.txt', '20208401000000-1200.txt', '20208501000000-1200.txt', '20208601000000-1200.txt', '20208701000000-1200.txt', '20208801000000-1200.txt', '20208901000000-1200.txt', '20209001000000-1200.txt', '20209101000000-1200.txt', '20209201000000-1200.txt', '20209301000000-1200.txt', '20209401000000-1200.txt', '20209501000000-1200.txt', '20209601000000-1200.txt', '20209701000000-1200.txt', '20209801000000-1200.txt', '20209901000000-1200.txt', '20210001000000-1200.txt', '20210101000000-1200.txt', '20210201000000-1200.txt', '20210301000000-1200.txt', '20210401000000-1200.txt', '20210501000000-1200.txt', '20210601000000-1200.txt', '20210701000000-1200.txt', '20210801000000-1200.txt', '20210901000000-1200.txt', '20211001000000-1200.txt', '20211101000000-1200.txt', '20211201000000-1200.txt', '20211301000000-1200.txt', '20211401000000-1200.txt', '20211501000000-1200.txt', '20211601000000-1200.txt', '20211701000000-1200.txt', '20211801000000-1200.txt', '20211901000000-1200.txt', '20212001000000-1200.txt', '20212101000000-1200.txt', '20212201000000-1200.txt', '20212301000000-1200.txt', '20212401000000-1200.txt', '20212501000000-1200.txt', '20212601000000-1200.txt', '20212701000000-1200.txt', '20212801000000-1200.txt', '20212901000000-1200.txt', '20213001000000-1200.txt', '20213101000000-1200.txt', '20213201000000-1200.txt', '20213301000000-1200.txt', '20213401000000-1200.txt', '20213501000000-1200.txt', '20213601000000-1200.txt', '20213701000000-1200.txt', '20213801000000-1200.txt', '20213901000000-1200.txt', '20214001000000-1200.txt', '20214101000000-1200.txt', '20214201000000-1200.txt', '20214301000000-1200.txt', '20214401000000-1200.txt', '20214501000000-1200.txt', '20214601000000-1200.txt', '20214701000000-1200.txt', '20214801000000-1200.txt', '20214901000000-1200.txt', '20215001000000-1200.txt', '20215101000000-1200.txt', '20215201000000-1200.txt', '20215301000000-1200.txt', '20215401000000-1200.txt', '20215501000000-1200.txt', '20215601000000-1200.txt', '20215701000000-1200.txt', '20215801000000-1200.txt', '20215901000000-1200.txt', '20216001000000-1200.txt', '20216101000000-1200.txt', '20216201000000-1200.txt', '20216301000000-1200.txt', '20216401000000-1200.txt', '20216501000000-1200.txt', '20216601000000-1200.txt', '20216701000000-1200.txt', '20216801000000-1200.txt', '20216901000000-1200.txt', '20217001000000-1200.txt', '20217101000000-1200.txt', '20217201000000-1200.txt', '20217301000000-1200.txt', '20217401000000-1200.txt', '20217501000000-1200.txt', '20217601000000-1200.txt', '20217701000000-1200.txt', '20217801000000-1200.txt', '20217901000000-1200.txt', '20218001000000-1200.txt', '20218101000000-1200.txt', '20218201000000-1200.txt', '20218301000000-1200.txt', '20218401000000-1200.txt', '20218501000000-1200.txt', '20218601000000-1200.txt', '20218701000000-1200.txt', '20218801000000-1200.txt', '20218901000000-1200.txt', '20219001000000-1200.txt', '20219101000000-1200.txt', '20219201000000-1200.txt', '20219301000000-1200.txt', '20219401000000-1200.txt', '20219501000000-1200.txt', '20219601000000-1200.txt', '20219701000000-1200.txt', '20219801000000-1200.txt', '20219901000000-1200.txt', '20220001000000-1200.txt', '20220101000000-1200.txt', '20220201000000-1200.txt', '20220301000000-1200.txt', '20220401000000-1200.txt', '20220501000000-1200.txt', '20220601000000-1200.txt', '20220701000000-1200.txt', '20220801000000-1200.txt', '20220901000000-1200.txt', '20221001000000-1200.txt', '20221101000000-1200.txt', '20221201000000-1200.txt', '20221301000000-1200.txt', '20221401000000-1200.txt', '20221501000000-1200.txt', '20221601000000-1200.txt', '20221701000000-1200.txt', '20221801000000-1200.txt', '20221901000000-1200.txt', '20222001000000-1200.txt', '20222101000000-1200.txt', '20222201000000-1200.txt', '20222301000000-1200.txt', '20222401000000-1200.txt', '20222501000000-1200.txt', '20222601000000-1200.txt', '20222701000000-1200.txt', '20222801000000-1200.txt', '20222901000000-1200.txt', '20223001000000-1200.txt', '20223101000000-1200.txt', '20223201000000-1200.txt', '20223301000000-1200.txt', '20223401000000-1200.txt', '20223501000000-1200.txt', '20223601000000-1200.txt', '20223701000000-1200.txt', '20223801000000-1200.txt', '20223901000000-1200.txt', '20224001000000-1200.txt', '20224101000000-1200.txt', '20224201000000-1200.txt', '20224301000000-1200.txt', '20224401000000-1200.txt', '20224501000000-1200.txt', '20224601000000-1200.txt', '20224701000000-1200.txt', '20224801000000-1200.txt', '20224901000000-1200.txt', '20225001000000-1200.txt', '20225101000000-1200.txt', '20225201000000-1200.txt', '20225301000000-1200.txt', '20225401000000-1200.txt', '20225501000000-1200.txt', '20225601000000-1200.txt', '20225701000000-1200.txt', '20225801000000-1200.txt', '20225901000000-1200.txt', '20226001000000-1200.txt', '20226101000000-1200.txt', '20226201000000-1200.txt', '20226301000000-1200.txt', '20226401000000-1200.txt', '20226501000000-1200.txt', '20226601000000-1200.txt', '20226701000000-1200.txt', '20226801000000-1200.txt', '20226901000000-1200.txt', '20227001000000-1200.txt', '20227101000000-1200.txt', '20227201000000-1200.txt', '20227301000000-1200.txt', '20227401000000-1200.txt', '20227501000000-1200.txt', '20227601000000-1200.txt', '20227701000000-1200.txt', '20227801000000-1200.txt', '20227901000000-1200.txt', '20228001000000-1200.txt', '20228101000000-1200.txt', '20228201000000-1200.txt', '20228301000000-1200.txt', '20228401000000-1200.txt', '20228501000000-1200.txt', '20228601000000-1200.txt', '20228701000000-1200.txt', '20228801000000-1200.txt', '20228901000000-1200.txt', '20229001000000-1200.txt', '20229101000000-1200.txt', '20229201000000-1200.txt', '20229301000000-1200.txt', '20229401000000-1200.txt', '20229501000000-1200.txt', '20229601000000-1200.txt', '20229701000000-1200.txt', '20229801000000-1200.txt', '20229901000000-1200.txt', '20230001000000-1200.txt', '20230101000000-1200.txt', '20230201000000-1200.txt', '20230301000000-1200.txt', '20230401000000-1200.txt', '20230501000000-1200.txt', '20230601000000-1200.txt', '20230701000000-1200.txt', '20230801000000-1200.txt', '20230901000000-1200.txt', '20231001000000-1200.txt', '20231101000000-1200.txt', '20231201000000-1200.txt', '20231301000000-1200.txt', '20231401000000-1200.txt', '20231501000000-1200.txt', '20231601000000-1200.txt', '20231701000000-1200.txt', '20231801000000-1200.txt', '20231901000000-1200.txt', '20232001000000-1200.txt', '20232101000000-1200.txt', '20232201000000-1200.txt', '20232301000000-1200.txt', '20232401000000-1200

```
In [21]: 1 #Open the child labour-related text files and show the contents
2 def process_files(file_paths, encoding='utf-8'):
3     for file_path in file_paths:
4         # Do something with each file
5         print("Processing file:", file_path)
6         with open(file_path, 'r', encoding=encoding, errors='ignore') as file:
7             text = file.read()
8             # Perform further processing or analysis on the file content
9             # For example, you can print the content of the file
10            print(text)
11            print("-" * 50) # Print a separator between files
12
13 # Call process_files with child_labor_folder list
14 process_files(child_labor_folder, encoding='utf-8')
15
```

Processing file: /Users/subaru/Library/CloudStorage/OneDrive-UniversityofLeeds/Courses/02\_Sem2/5403\_Creative Coding/Week4/Alex/dataset/20151201211500-2747.txt

woodrow wilson wasn't morally or ethically pure. no president is.

bangor, me

33

clouds

close

login  
logout  
subscribe

search for:

search

```
In [22]: 1 import pandas as pd
2 def read_text_files(file_paths):
3     data = []
4     for file_path in file_paths:
5         with open(file_path, 'r', encoding='utf-8') as file:
6             content = file.read()
7             data.append(content)
8     return data
9
10 # Assign a dataframe and accommodate contents.
11 def text_files_to_dataframe(file_paths):
12     # Read text files
13     texts = read_text_files(file_paths)
14
15     # Create DataFrame
16     df = pd.DataFrame({'Text': texts})
17     # Optionally, add a column for file names
18     df['File_Name'] = [os.path.basename(file_path) for file_path in file_paths]
19
20
21
22 # Create DataFrame
23 df = text_files_to_dataframe(child_labor_folder)
24
25 # Display DataFrame
26 print(df)
27
```

Text \nwoodrow wilson wasn't morally or ethically p...  
1 hp is turning trash into printer cartridgeseng...  
2 npakistan's censor has a thing or two to say ...  
3 president obama talks about human rights, bett...  
4 what we uncovered about ivanka trump's fashion...  
...  
924 migrant construction workers' children in thai...  
925 top aid groups call on obama to help refugeesx...  
926 \nis your work killing you? \nnskip to the co...  
927 \nyou can't have financial inclusion without d...  
928 \nhershey to invest \$0.5bn by 2030 for sustain...

File\_Name  
0 20151201211500-2747.txt  
1 20170615204500-612.txt  
2 20160609040000-247.txt  
3 20160525113000-2449.txt  
4 20170209214500-1865.txt  
...  
924 20180329113000-320.txt  
925 20160920184500-1590.txt  
926 20180424141500-2753.txt  
927 20160715114500-1237.txt  
928 20180405150000-996.txt

[929 rows x 2 columns]

```
In [244]: 1 ## Search for specific articles
2 test_keywords = 'United States'
3 test = df[df['Text'].str.contains(test_keywords, case=False)]
4 test
```

	Text	File_Name
0	\nwoodrow wilson wasn't morally or ethically p...	20151201211500-2747.txt
3	president obama talks about human rights, bett...	20160525113000-2449.txt
5	\nthai pm halts migrant labor law after protes...	20170705150000-2474.txt
7	\neducation attainment: muslims, hindus lag fa...	20161214084500-1377.txt
8	\na colorful bus brings books and joy to afgha...	20180316160000-2105.txt
...	...	...
916	top 8 countries with most horrific human right...	20180423150000-1619.txt
919	\nhuman rights archives - u.s. embassy in tuni...	20170612120000-2284.txt
922	betsy devos: religion and profit in the war on...	20161213081500-1022.txt
923	\nstudents pulled between work and school in p...	20161130143000-972.txt
926	\nis your work killing you? \nnskip to the co...	20180424141500-2753.txt

301 rows x 2 columns

```
In [254]: 1 from IPython.core.interactiveshell import InteractiveShell
2 InteractiveShell.ast_node_interactivity = "all"
```

```
In [265]: 1 df[df.index==65]
```

```
Out[265]:
```

Text	File_Name
65 \na loss of nerve   commonweal magazine\n\nmail... 20160117094500-947.txt	

```
In [31]: 1 df.iloc[100]
```

```
Out[31]: "guatemala rescues 22 children from forced labor: police | reuters\n\nskip to main contentexclusive news, data and analytics for financial market professionalslearn more aboutreuters homepagebusinessmarketsustainabilitylegalbreakingviewstechnologyinvestigationsmoremy viewregisterworldguatemala rescues 22 children from forced labor: policeby anastasia moloneyjune 1, 20174:59 pm utcupdated ago bogota (t homson reuters foundation) - guatemalan police have rescued 22 children - some working 16-hour days making corn tortillas - as part of a crackdown on widespread child labor in the central american nation, accord ing to public prosecutors.police arrested six people on charges of child labor on wednesday, following 37 raids on shops and food stalls in the capital guatemala city.the children rescued were aged 14 and 17, and mo st came from impoverished indigenous communities in rural areas.children were forced to work seven days a week in unhealthy conditions, and with little pay and food, prosecutors said.the police raids aimed to targ et "those whose take advantage of a child's vulnerability to get cheap labor," the public prosecutor's office said in a statement on wednesday.poverty drives child labor in guatemala.about half of guatemala's 15 million people live in poverty and the country has one of the world's highest rates of chronic malnutrition, mostly affecting the children of indigenous communities.to put food on the table, some parents encourage their children to quit school and work in the capital, leaving them open to sexual exploitation and forced labor, experts say.nearly 415,000 children, some as young as seven, work in guatemala, according to the u.s.government.most child labor exists in agriculture, along with domestic work, forced begging and garbage scavenging.it is common to see children working day and night as shoe shiners and street vendors on the cap ital's streets.worldwide there are 168 million children in child labor, according to the international labour organization.the united nations has a global goal to end all child labor by 2025 and to eradicate forced labor and slavery by 2030.reporting by anastasia moloney @anastasiabogota, editing by lyndsay griffiths. please credit the thomson reuters foundation, the charitable arm of thomson reuters, that covers humanitarian n ews, women's rights, trafficking, property rights, climate change and resilience. visit <a href="http://news.trust.org">news.trust.org</a>our standards: the thomson reuters trust principles. acquire licensing right s, opens new tabread nextunitexted statescategoryinfluential koch group endorses haley's 2024 republican presidential bid6:53 pm utc article with videoeuropecategoryukraine says wife of spymaster budanov was po isoned7:03 pm utc · updated ago worldcategoryiran finalises deal to buy russian fighter jets - tasnim9:04 am utc article with galleryeuropacategorycourt says public employees may be barred from wearing head s carf1:35 pm utc more from reutersworldweek after shock wilders win, dutch government talks set to begineuropacategory · november 28, 2023 · 11:53 am utc talks to form a new dutch government were set to start on tuesday, almost a week after the upset election victory of anti-islam populist geert wilders, with a former government minister picked to sound out workable coalitions.article with galleryeuropacategorygreece, britain trade blame over cancelled meeting in parthenon marbles dispute5:49 pm utc article with galleryeuropacategorygerman authorities urge people to stay home amid deadly winter weather10:51 am utc unitied s tatescategoryembattled republican george santos faces fresh move to oust him from house7:15 pm utc · updated ago article with videoworldcategoryisrael-hamas war: answering readers' questions4:59 pm utc site indexboxeworldbusinessmarketssustainabilitylegalbreakingviewstechnologyinvestigations, opens new tabsportssciencelifestyleabout reutersabout reuters, opens new tabcareers, opens new tabreuters news agen cy, opens new tabbrand attribution guidelines, opens new tabreuters leadership, opens new tabreuters fact check, opens new tabreuters diversity report, opens new tabstay informedownload the app (ios), opens ne w tabdownload the app (android), opens new tabnewsletters, opens new tabinformation you can trustreuters, the news and media division of thomson reuters, is the world's largest multimedia news provider, reac hing billions of people worldwide every day, reuters provides business, financial, national and international news to professionals via desktop terminals, the world's media organizations, industry events and directly t o consumers.follow usthomson reuters productswestlaw, opens new tabbuild the strongest argument relying on authoritative content, attorney-editor expertise, and industry defining technology.onsource, opens ne w tabthe most comprehensive solution to manage all your complex and ever-expanding tax and compliance needs.checkpoint, opens new tabthe industry leader for online information for tax, accounting and financ e professionals.leg productsworkspace, opens new tab access unmatched financial data, news and content in a highly-customised workflow experience on desktop, web and mobile.data catalogue, opens new tab br owsse an unrivaled portfolio of real-time and historical market data and insights from worldwide sources and experts.world-check, opens new tabscreen for heightened risk individual and entities globally to help unco ver hidden risks in business relationships and human networks.advertise with us, opens new tabadvertising guidelines, opens new tabcoupons, opens new tabacquire licensing rights, opens new tabcookies, opens ne w tabterms of use, opens new tabprivacy, opens new tabdigital accessibility, opens new tabcorrections, opens new tabsite feedback, opens new taball quotes delayed a minimum of 15 minutes. see here for a compre hensive list of exchanges and delays.© 2023 reuters. all rights reserved"
```

```
In [30]: 1 df_text_file = df['File_Name']  
2 df_text_file.to_csv('df_text_file.csv', index=False, encoding='utf-8')
```

## Step3: Analysis with ChatGPT

```
In [16]: 1 import os  
2 import openai  
3 from openai import OpenAI  
4 from openai.types import Completion, CompletionChoice, CompletionUsage  
5  
6 API_KEY= "XXXXX"
```

### 3-1: Create function to identify Companies

```
In [72]: 1 def chat_gpt_companies(text, API_KEY):  
2     client = OpenAI(api_key = API_KEY)  
3     response = client.chat.completions.create(  
4         messages=[  
5             {  
6                 "role": "system",  
7                 "content": f"""Can you find the companies that are mentioned in relation to child labour exploitation?  
8                 Here is the text:{text}。  
9                 Please provide a list of company names, without any additional text or sentences.  
10                If there are no relevant company names found, please return 'None'.  
11                Here is an example of your message: 'Microsoft, Alphabet'.  
12                """  
13            },  
14            {  
15                "role": "user",  
16                "content": text  
17            }  
18        ],  
19        model = "gpt-3.5-turbo"  
20    )  
21  
22    return (response.choices[0].message.content)
```

### 3-2: Create function to identify Commodities

```
In [73]: 1 ## function to extract commodity info using chat gpt api  
2 def chat_gpt_commodities(text, API_KEY):  
3     client = OpenAI(api_key = API_KEY)  
4     response = client.chat.completions.create(  
5         messages=[  
6             {  
7                 "role": "system",  
8                 "content": f"""Can you find the commodities that are mentioned in relation to child labour exploitation?  
9                 Here is the text:{text}。  
10                Please provide a list of commodity names, without any additional text or sentences, using only the names in the commodity list: '{commodity_labour}'.  
11                If there are no relevant commodity names found, please return 'None'.  
12                Here is an example of your message: 'Jewelry, Apparel'.  
13                """  
14            },  
15            {  
16                "role": "user",  
17                "content": text  
18            }  
19        ],  
20        model = "gpt-3.5-turbo"  
21    )  
22  
23    return (response.choices[0].message.content)
```

### 3-3: Create function to identify Countries

```
In [74]: 1 ## function to extract country info using chat gpt api
2 def chat_gpt_countries(text, API_KEY):
3     client = OpenAI(api_key = API_KEY)
4     response = client.chat.completions.create(
5         messages=[
6             {
7                 "role": "system",
8                 "content": f"""Can you find the countries that are mentioned in relation to child labour exploitation?
9                 Here is the text: '{text}'.
10                Please provide a list of country names, without any additional text or sentences.
11                If there are no relevant country names found, please return 'None'.
12                Here is an example of your message: 'United States, Indonesia'.
13                """
14            },
15            {
16                "role": "user",
17                "content": text
18            }
19        ],
20        model = "gpt-3.5-turbo"
21    )
22
23    return (response.choices[0].message.content)
24
```

### 3-4: Run all the ChatGPT functions and Store returned result into DataFrame

CAUTION!!!: Please DO NOT RUN the following chunk as it consumes a lot of tokens.

```
In [83]: 1 # CAUTION! Please DO NOT RUN this step as it consumes a lot of tokens#####
2 #####
3
4 import pandas as pd
5
6 # Assuming you have created an empty DataFrame named output_df
7 #output_df = pd.DataFrame(columns=['Companies', 'Commodities', 'Countries', 'Additional_Info'])
8 output_df = pd.DataFrame(columns=['Companies', 'Commodities', 'Countries'])
9
10 # Assuming you have a DataFrame named input_df with a column named 'text' containing the text data
11 start_index = 0
12 for index, row in df.iloc[start_index: ].iterrows():
13     text = row["text"]
14     try:
15         output1 = chat_gpt_companies(text, API_KEY)
16         output2 = chat_gpt_commodities(text, API_KEY)
17         output3 = chat_gpt_countries(text, API_KEY)
18         print("Index: {}, Companies: {} Commodities: {} Countries: {}".format(index, output1, output2, output3))
19         output_df.loc[index] = [output1, output2, output3]
20     except Exception as e:
21         print(f"Error processing text index: {index}.")
22         continue
23
24
25 # Now output_df contains the combined outputs
26 print(output_df)
27
```

Index: 0, Companies: None Commodities: None Countries: None  
Index: 1, Companies: None Commodities: Recycled materials Countries: Haiti, Honduras  
Index: 2, Companies: None Commodities: Children, Child Labor, Child Abuse, Contraceptives. Countries: Pakistan  
Index: 3, Companies: None Commodities: None Countries: United States, Vietnam, China  
Index: 4, Companies: G-III Apparel Group Commodities: Apparel, Textiles, Garments, Fashion Countries: Indonesia, China, Vietnam  
Index: 5, Companies: None Commodities: Fisheries, Construction, Agriculture Countries: Myanmar, Laos, Cambodia  
Index: 6, Companies: None Commodities: Textiles, Apparel Countries: Kyrgyzstan  
Index: 7, Companies: None Commodities: None Countries: United States, Israel, India, Nepal, Bangladesh, Sub-Saharan Africa, Europe, Middle East  
Index: 8, Companies: None Commodities: Books Countries: Afghanistan  
Error processing text index: 9.  
Index: 10, Companies: None Commodities: Corn, Food Processing Products, Agricultural Produce Countries: Guatemala, United States  
Index: 11, Companies: None Commodities: 'Child labor' Countries: Mexico, United States  
Index: 12, Companies: None Commodities: None Countries: Pakistan  
Index: 13, Companies: None Commodities: None Countries: United States, Central America  
Index: 14, Companies: Nike, Gap, Commodities: Cotton, Textiles, Apparel, Garments Countries: Bangladesh, Uzbekistan, India  
Index: 15, Companies: None Commodities: Ores, Gold, Textiles, Apparel, Electronics Countries: Brazil, Canada  
Index: 16, Companies: Disney Commodities: Newspaper Countries: United States, New York  
Index: 17, Companies: None Commodities: None Countries: Afghanistan, Pakistan  
Index: 18, Companies: Thai Union, Thai Frozen Foods Association Commodities: Shrimp, Seafood, Migrant laborers Countries: Thailand, Myanmar, United States, Europe, Australia, Indonesia

### 3-5: Add summary column to DataFrame that contains involved companies, commodities and countries in each case

```
In [148]: 1 output_df['comb'] = "Case" + output_df.index.astype(str) + ";<<" + 'Companies: ' + output_df['Companies'].astype(str) + '. Commodities: ' + output_df['Commodities'].astype(str) + '. Countries: ' + output_
2 output_df
```

Out[148]:

	Companies	Commodities	Countries	comb
0	None	None		None "Case0":<<Companies: None. Commodities: None. ...
1	None	Recycled materials	Haiti, Honduras	"Case1":<<Companies: None. Commodities: Recycl...
2	None	Children, Child Labor, Child Abuse, Contracept...	Pakistan	"Case2":<<Companies: None. Commodities: Childr...
3	None	None	United States, Vietnam, China	"Case3":<<Companies: None. Commodities: None. ...
4	G-III Apparel Group	Apparel, Textiles, Garments, Fashion	Indonesia, China, Vietnam	"Case4":<<Companies: G-III Apparel Group. Comm...
...	...	...	...	...
924	None	None	Cambodia, Myanmar, Laos	"Case924":<<Companies: None. Commodities: None...
925	None.	None.	Syria, South Sudan, Yemen, Burundi, Central Af...	"Case925":<<Companies: None.. Commodities: Non...
926	None	None	United States, United Kingdom	"Case926":<<Companies: None. Commodities: None...
927	None	None	None	"Case927":<<Companies: None. Commodities: None...
928	hershey	Cocoa, Food processing products, Agricultural ...	Côte d'Ivoire, Ghana	"Case928":<<Companies: hershey. Commodities: C...

908 rows × 4 columns

### 3-6: Export the result as CSV file.

```
In [149]: 1 output_df.to_csv('output_data2.csv', index=False)
```

In [150]: 1 output\_df

	Companies	Commodities	Countries	comb
0	None	None	None	"Case0":<<Companies: None. Commodities: None. ...
1	None	Recycled materials	Haiti, Honduras	"Case1":<<Companies: None. Commodities: Recycl...
2	None	Children, Child Labor, Child Abuse, Contracept...	Pakistan	"Case2":<<Companies: None. Commodities: Childr...
3	None	None	United States, Vietnam, China	"Case3":<<Companies: None. Commodities: None. ...
4	G-III Apparel Group	Apparel, Textiles, Garments, Fashion	Indonesia, China, Vietnam	"Case4":<<Companies: G-III Apparel Group. Comm...
...	...	...	...	...
924	None	None	Cambodia, Myanmar, Laos	"Case924":<<Companies: None. Commodities: None...
925	None.	None.	Syria, South Sudan, Yemen, Burundi, Central Af...	"Case925":<<Companies: None.. Commodities: Non...
926	None	None	United States, United Kingdom	"Case926":<<Companies: None. Commodities: None...
927	None	None	None	"Case927":<<Companies: None. Commodities: None...
928	hershey	Cocoa, Food processing products, Agricultural ...	Côte d'Ivoire, Ghana	"Case928":<<Companies: hershey. Commodities: C...

908 rows × 4 columns

In [151]: 1 output\_df['comb'][30]

Out[151]: "Case30":<<Companies: None. Commodities: cocoa. Countries: Ivory Coast>> <br>"

## Step4: Summary Analysis with ChatGPT

### 4-1: Data preparation for feeding data into ChatGPT

```
In [19]: 1 # Extracting column values into a list
2 list_companies = output_df['Companies'].tolist()
3 string_companies = ', '.join(list_companies)
4
5 # Extracting column values into a list
6 list_commodities = output_df['Commodities'].tolist()
7 string_commodities = ', '.join(list_commodities)
8
9 # Extracting column values into a list
10 list_countries = output_df['Countries'].tolist()
11 string_countries = ', '.join(list_countries)
12
13 ## Extracting column values into a list
14 list_comb = output_df['comb'].tolist()
15 string_comb = ', '.join(list_comb)
16 print(string_comb)
17
```

NameError Traceback (most recent call last)

Cell In[19], line 4  
1 # !CAUTION! Please DO NOT RUN this step as it consumes a lot of tokens.#####  
2  
3 # Extracting column values into a list  
4 list\_companies = output\_df['Companies'].tolist()  
5 string\_companies = ', '.join(list\_companies)  
7 # Extracting column values into a list

NameError: name 'output\_df' is not defined

```
In [154]: 1 string_companies
```

## 4-2: Summary analysis for top 20 companies with ChatGPT

```
In [156]: 1 #Companies
2 client = OpenAI(api_key = API_KEY)
3 response = client.chat.completions.create(
4     messages=[

5         {
6             "role": "system",
7             "content": """
8                 Your role is to summarise the list of companies involved in child labour.
9                 Here is the list: """ + string_companies
10            },
11        ],
12        [
13            {
14                "role": "user",
15                "content": """
16                    Given the list provided, please rank the top 20 companies that are most likely to be found in relation to child exploitation.
17                    Please ignore 'None' in the list. """
18            }
19        ],
20        model = "gpt-3.5-turbo"
21    )
22 ranked_companies = response.choices[0].message.content
23 print(ranked_companies)
```

Based on the list provided, the top 20 companies most likely to be found in relation to child exploitation are as follows:

- 1. Nestle
  - 2. Cargill
  - 3. Walmart
  - 4. Nike
  - 5. Apple
  - 6. Hershey
  - 7. Mars
  - 8. Ferrero
  - 9. Driscoll's
  - 10. H&M
  - 11. Inditex
  - 12. Gap
  - 13. Pepsico
  - 14. Procter & Gamble
  - 15. McDonald's
  - 16. Coca-Cola
  - 17. Samsung
  - 18. Amazon
  - 19. Starbucks
  - 20. Disney

These companies have appeared most frequently in the list in relation to child exploitation.

#### 4-3: Summary analysis for top 20 commodities with ChatGPT

```
In [157]: 1 #Commodities
2 client = OpenAI(api_key = API_KEY)
3 response = client.chat.completions.create(
4     messages=[
5
6         {
7             "role": "system",
8             "content": """
9             Your role is to summarise the list of commodities involved in child labour.
10            Here is the list: """ + string_commodities
11        },
12        {
13            "role": "user",
14            "content": """
15            Given the list provided, please rank the top 20 commodities that are most likely to be found in relation to child exploitation.
16            Please ignore 'None' in the list. """
17        }
18    ],
19    model = "gpt-3.5-turbo"
20 )
21
22 ranked_commodities = response.choices[0].message.content
23 print(ranked_commodities)
```

Based on the list of commodities mentioned in relation to child labor exploitation, the top 20 commodities likely to be found in such situations are:

1. Textiles
2. Apparel
3. Footwear
4. Toys
5. Handicrafts
6. Electronics
7. Garments
8. Furniture
9. Ores
10. Bricks
11. Carpentry
12. Carpets
13. Paper products
14. Agricultural produce
15. Medical supplies
16. Packaging materials
17. Metal products
18. Building materials
19. Leather goods
20. Cosmetics

#### 4-4: Summary analysis for top 20 countries with ChatGPT

```
In [158]: 1 #Countries
2 client = OpenAI(api_key = API_KEY)
3 response = client.chat.completions.create(
4     messages=[
5
6         {
7             "role": "system",
8             "content": """
9             Your role is to summarise the list of countries involved in child labour.
10            Here is the list: """ + string_countries
11        },
12        {
13            "role": "user",
14            "content": """
15            Given the list provided, please rank the top 20 countries that are most likely to be found in relation to child exploitation.
16            Please ignore 'None' in the list. """
17        }
18    ],
19    model = "gpt-3.5-turbo"
20 )
21
22 ranked_countries = response.choices[0].message.content
23 print(ranked_countries)
```

Based on the extensive list provided, the top 20 countries that are most likely to be found in relation to child exploitation are:

1. India
2. United States
3. Pakistan
4. Bangladesh
5. Syria
6. Afghanistan
7. Iraq
8. Nigeria
9. Mexico
10. Indonesia
11. Democratic Republic of Congo
12. Vietnam
13. Yemen
14. Thailand
15. Russia
16. China
17. Lebanon
18. Sudan
19. Mali
20. Ghana

#### 4-5: Export the results

```
In [159]: 1 # File path
2 file_path = "ranked_output2.txt"
3
4 # Open the file in write mode
5 with open(file_path, "w") as file:
6     # Write the strings to the file
7     file.write("Worst Countries are" + ranked_countries + "\n")
8     file.write("Most common companies are" + ranked_companies + "\n")
9     file.write("Most common commodities are" + ranked_commodities + "\n")
10
11 print("File created successfully.")
12
```

File created successfully.

## Step5: Visualisation

```
In [53]: 1 output_df = pd.read_csv("output_data2.csv")
2 output_df
```

Out[53]:

	Companies	Commodities	Countries	comb
0	None	None	None	"Case0":<<Companies: None. Commodities: None. ...
1	None	Recycled materials	Haiti, Honduras	"Case1":<<Companies: None. Commodities: Recycl...
2	None	Children, Child Labor, Child Abuse, Contracept...	Pakistan	"Case2":<<Companies: None. Commodities: Childr...
3	None	None	United States, Vietnam, China	"Case3":<<Companies: None. Commodities: None. ...
4	G-III Apparel Group	Apparel, Textiles, Garments, Fashion	Indonesia, China, Vietnam	"Case4":<<Companies: G-III Apparel Group. Comm...
...	...	...	...	...
903	None	None	Cambodia, Myanmar, Laos	"Case924":<<Companies: None. Commodities: None...
904	None.	None.	Syria, South Sudan, Yemen, Burundi, Central Af...	"Case925":<<Companies: None.. Commodities: Non...
905	None	None	United States, United Kingdom	"Case926":<<Companies: None. Commodities: None...
906	None	None	None	"Case927":<<Companies: None. Commodities: None...
907	hershey	Cocoa, Food processing products, Agricultural ...	Côte d'Ivoire, Ghana	"Case928":<<Companies: hershey. Commodities: C...

908 rows × 4 columns

## 5-1: Countries

```
In [126]: 1 output_df['Countries'] = output_df['Countries'].str.replace('.','')
2 output_df
```

/var/folders/ct/fb9zr96x3g1dxcqhdtydzv5h0000gn/T/ipykernel\_3558/3623789606.py:1: FutureWarning:

The default value of regex will change from True to False in a future version. In addition, single character regular expressions will \*not\* be treated as literal strings when regex=True.

Out[126]:

	Companies	Commodities	Countries	comb
0	None	None	None	"Case0":<<Companies: None. Commodities: None. ...
1	None	Recycled materials	Haiti, Honduras	"Case1":<<Companies: None. Commodities: Recycl...
2	None	Children, Child Labor, Child Abuse, Contracept...	Pakistan	"Case2":<<Companies: None. Commodities: Childr...
3	None	None	United States, Vietnam, China	"Case3":<<Companies: None. Commodities: None. ...
4	G-III Apparel Group	Apparel, Textiles, Garments, Fashion	Indonesia, China, Vietnam	"Case4":<<Companies: G-III Apparel Group. Comm...
...	...	...	...	...
903	None	None	Cambodia, Myanmar, Laos	"Case924":<<Companies: None. Commodities: None...
904	None.	None.	Syria, South Sudan, Yemen, Burundi, Central Af...	"Case925":<<Companies: None.. Commodities: Non...
905	None	None	United States, United Kingdom	"Case926":<<Companies: None. Commodities: None...
906	None	None	None	"Case927":<<Companies: None. Commodities: None...
907	hershey	Cocoa, Food processing products, Agricultural ...	Côte d'Ivoire, Ghana	"Case928":<<Companies: hershey. Commodities: C...

908 rows × 4 columns

```
In [259]: 1 output_df[output_df.index==13]
```

Out[259]:

	Companies	Commodities	Countries	comb
13	Nike, Gap	Cotton, Textiles, Apparel, Garments	Bangladesh, Uzbekistan, India	"Case14":<<Companies: Nike, Gap.. Commodities:...

```
In [128]: 1 # test_country = ['Congo']
2 # test = output_df[output_df['Countries'].isin(test_country)]
3 # test
```

```
In [135]: 1 # Unify the country names
2
3 ## DRC
4 output_df['Countries'] = output_df['Countries'].str.replace('Democratic Republic of Congo (DRC)', 'Democratic Republic of Congo')
5 output_df['Countries'] = output_df['Countries'].str.replace('Democratic Republic of the (DRC)', 'Democratic Republic of Congo')
6 output_df['Countries'] = output_df['Countries'].str.replace('Democratic Republic of the Congo', 'Democratic Republic of Congo')
7 output_df['Countries'] = output_df['Countries'].str.replace('DR Congo', 'Democratic Republic of Congo')
8
9 # Germany
10 output_df['Countries'] = output_df['Countries'].str.replace('Germany', 'Germany')
11 # Nicaragua
12 output_df['Countries'] = output_df['Countries'].str.replace('Nicaragua', 'Nicaragua')
13
14 # United States
15 output_df['Countries'] = output_df['Countries'].str.replace('America', 'United States')
16
17 # Brunei
18 output_df['Countries'] = output_df['Countries'].str.replace('Brunei', 'Brunei Darussalam')
19 output_df['Countries'] = output_df['Countries'].str.replace('and Brunei Darussalam', 'Brunei Darussalam')
20 output_df['Countries'] = output_df['Countries'].str.replace('Brunei Darussalam Darussalam', 'Brunei Darussalam')
21
22
23 # North Korea
24 output_df['Countries'] = output_df['Countries'].str.replace('North Korea', 'North Korea')
25
26 # United Kingdom
27 output_df['Countries'] = output_df['Countries'].str.replace('Britain', 'United Kingdom')
28 output_df['Countries'] = output_df['Countries'].str.replace('United Kingdom (UK)', 'United Kingdom')
29 output_df['Countries'] = output_df['Countries'].str.replace('the UK', 'United Kingdom')
30 output_df['Countries'] = output_df['Countries'].str.replace('UK', 'United Kingdom')
31
32 # Venezuela
33 output_df['Countries'] = output_df['Countries'].str.replace('Venezuela (Bolivarian Republic of)', 'Venezuela')
34
35 # Myanmar
36 output_df['Countries'] = output_df['Countries'].str.replace('Myanmar (Burma)', 'Myanmar')
37
38 # Russia
39 output_df['Countries'] = output_df['Countries'].str.replace('Russian Federation', 'Russia')
40
41
```

/var/folders/ct/fb9zr96x3g1dxcqhdtdzv5h0000gn/T/ipykernel\_3558/2142745092.py:4: FutureWarning:

The default value of regex will change from True to False in a future version.

/var/folders/ct/fb9zr96x3g1dxcqhdtdzv5h0000gn/T/ipykernel\_3558/2142745092.py:5: FutureWarning:

The default value of regex will change from True to False in a future version.

/var/folders/ct/fb9zr96x3g1dxcqhdtdzv5h0000gn/T/ipykernel\_3558/2142745092.py:28: FutureWarning:

The default value of regex will change from True to False in a future version.

/var/folders/ct/fb9zr96x3g1dxcqhdtdzv5h0000gn/T/ipykernel\_3558/2142745092.py:33: FutureWarning:

The default value of regex will change from True to False in a future version.

/var/folders/ct/fb9zr96x3g1dxcqhdtdzv5h0000gn/T/ipykernel\_3558/2142745092.py:36: FutureWarning:

The default value of regex will change from True to False in a future version.

In [136]: 1 output\_df['Countries'].value\_counts()

```
Out[136]: None 166
United States 65
Pakistan 18
India 17
Democratic Republic of Congo 13
Turkey, Syria, Germany, United Kingdom ...
Qatar, United States, Thailand, Philippines, Brazil, Canada, United Kingdom 1
United States, Bosnia, Russia, Norway, France, United Kingdom, Portugal, Spain, Japan, Lithuania, Greece, Niger, Angola, Mali, Central African Republic, Somalia 1
Mexico, United States, Arizona, California, Texas, Washington, Florida, Oregon, North Carolina 1
United States, United Kingdom 1
Name: Countries, Length: 476, dtype: int64
```

```
In [137]: 1 country_counts={}
2 for data in output_df['Countries']:
3     data = data.split(',')
4     for country in data:
5         if country in country_counts:
6             country_counts[country] += 1
7         else:
8             country_counts[country] = 1
9
```

In [138]: 1 country\_counts

```
Out[138]: {'None': 199,
'Haiti': 14,
'Honduras': 10,
'Pakistan': 63,
'United States': 243,
'Vietnam': 38,
'China': 86,
'Indonesia': 38,
'Myanmar': 30,
'Laos': 7,
'Cambodia': 18,
'Kyrgyzstan': 3,
'Israel': 14,
'India': 97,
'Nepal': 13,
'Bangladesh': 54,
'Sub-Saharan Africa': 2,
'Europe': 5,
'Middle East': 1,
```

```
In [139]: 1 country_counts
```

```
Out[139]: {'None': 199,
'Haiti': 14,
'Honduras': 10,
'Pakistan': 63,
'United States': 243,
'Vietnam': 38,
'China': 86,
'Indonesia': 38,
'Myanmar': 30,
'Laos': 7,
'Cambodia': 18,
'Kyrgyzstan': 3,
'Israel': 14,
'India': 97,
'Nepal': 13,
'Bangladesh': 54,
'Sub-Saharan Africa': 2,
'Europe': 5,
'Middle East': 1,
'...': 120}
```

```
In [140]: 1 country_counts['Democratic Republic of Congo (DRC)']
```

```
Out[140]: 2
```

```
In [141]: 1 country_counts_df = pd.DataFrame(list(country_counts.items()), columns=['Country', 'Count'])
2 country_counts_df = country_counts_df.sort_values(by='Count', ascending=False)
3
4 # Remove "None"
5 country_counts_df = country_counts_df[country_counts_df['Country'] != 'None']
6 country_counts_df
```

```
Out[141]:
```

	Country	Count
4	United States	243
13	India	97
6	China	86
45	Syria	64
3	Pakistan	63
...	...	...
287	Ethiopia\n	1
286	Hague Conference	1
285	United Nations (UN)	1
170	Venezuela (Bolivarian Republic of)	1
371	Former Soviet States	1

371 rows × 2 columns

```
In [142]: 1 country_counts_df.to_csv('country_counts.csv', index=False, encoding='utf-8')
```

```
In [143]: 1 country_counts_df_top20 = country_counts_df[country_counts_df['Count']>=20]
2 country_counts_df_top20
```

```
Out[143]:
```

	Country	Count
4	United States	243
13	India	97
6	China	86
45	Syria	64
3	Pakistan	63
21	Mexico	58
15	Bangladesh	54
19	Afghanistan	39
55	Ghana	39
52	Russia	39
7	Indonesia	38
51	United Kingdom	38
5	Vietnam	38
72	Turkey	37
47	Iraq	36
23	Uzbekistan	35
34	Democratic Republic of Congo	33
8	Myanmar	30
48	Jordan	30
69	Malaysia	29
49	Lebanon	29
37	Japan	27
25	Canada	26
63	Nigeria	25
27	Thailand	25
86	Peru	21
28	Australia	21
113	Philippines	21
97	Sierra Leone	20
65	Iran	20

```
In [144]: 1 fig = px.choropleth(country_counts_df_top20,
2                         locations='Country',
3                         locationmode='country names',
4                         color='Count',
5                         hover_name='Country',
6                         color_continuous_scale='Reds',
7                         title='Frequencies of child exploitation news by country (over 20 counts only)')
8 fig.update_layout(width=800, height=500)
9 fig.write_html("country_counts.html")
10 fig.show()
```

## 5-2: Commodities

```
In [200]: 1 # Remove unnecessary punctuation marks:
2 output_dff['Commodities'] = output_dff['Commodities'].str.replace("\n", ",")
3 output_dff['Commodities'] = output_dff['Commodities'].str.replace("[", ",")
4 output_dff['Commodities'] = output_dff['Commodities'].str.replace("]", ",")
5 output_dff['Commodities'] = output_dff['Commodities'].str.replace(".", ",")
6 output_dff['Commodities'] = output_dff['Commodities'].str.replace("", ",")
7
8 # Unify commodity names
9 ## Cobalt
10 output_dff['Commodities'] = output_dff['Commodities'].str.replace("cobolt", "cobalt")
11 ## Apparel
12 output_dff['Commodities'] = output_dff['Commodities'].str.replace("cobolt", "cobalt")
13
14
15
```

/var/folders/ct/fb9zr96x3g1dxcqhdtdzv5h0000gn/T/ipykernel\_3558/2738018453.py:3: FutureWarning:

The default value of regex will change from True to False in a future version. In addition, single character regular expressions will \*not\* be treated as literal strings when regex=True.

/var/folders/ct/fb9zr96x3g1dxcqhdtdzv5h0000gn/T/ipykernel\_3558/2738018453.py:4: FutureWarning:

The default value of regex will change from True to False in a future version. In addition, single character regular expressions will \*not\* be treated as literal strings when regex=True.

/var/folders/ct/fb9zr96x3g1dxcqhdtdzv5h0000gn/T/ipykernel\_3558/2738018453.py:5: FutureWarning:

The default value of regex will change from True to False in a future version. In addition, single character regular expressions will \*not\* be treated as literal strings when regex=True.

```
In [201]: 1 test_country = ['Cobolt']
2 test = output_df[output_dff['Commodities'].isin(test_country)]
3 test
```

Out[201]:

	Companies	Commodities	Countries	comb
270		None	Cobolt	Democratic Republic of Congo "Case279":<<Companies: None. Commodities: Cobo...
565	Samsung, Huawei, Microsoft.	Cobolt		Congo "Case579":<<Companies: Samsung, Huawei, Micros...
667		None	Cobolt	Congo "Case684":<<Companies: None. Commodities: ['Co...
781	GM, Renault-Nissan, Fiat Chrysler, Volkswagen,...	Cobolt	Democratic Republic of Congo (DRC)	"Case800":<<Companies: GM, Renault-Nissan, Fia...
845	Apple, HP, Samsung SDI, Sony, Dell, Foxconn, F...	Cobolt		Congo "Case865":<<Companies: Apple, HP, Samsung SDI,...

```
In [202]: 1 commodity_counts = {}
2 for data in output_dff['Commodities']:
3     data = data.split(',')
4     for commodity in data:
5         if commodity in commodity_counts:
6             commodity_counts[commodity] += 1
7         else:
8             commodity_counts[commodity] = 1
```

In [203]: 1 commodity\_counts

```
Out[203]: {'None': 465,
'Recycled materials': 8,
'Children': 3,
'Child Labor': 10,
'Child Abuse': 2,
'Contraceptives': 1,
'Apparel': 115,
'Textiles': 101,
'Garments': 61,
'Fashion': 1,
'Fisheries': 2,
'Construction': 4,
'Agriculture': 6,
'Books': 3,
'Corn': 4,
'Food Processing Products': 1,
'Agricultural Produce': 3,
'Child labor': 31,
'Cotton': 37,
'...': 32}
```

```
In [204]: 1 commodity_counts_df = pd.DataFrame(list(commodity_counts.items()), columns=['Commodity', 'Count'])
2
3 # Remove "None" and inappropriate words
4 commodity_counts_df = commodity_counts_df[~commodity_counts_df['Commodity'].isin(['None', 'Child labor', 'Child Labor'])]
5 # Reorganise
6 commodity_counts_df = commodity_counts_df.sort_values(by='Count', ascending=False)
7 commodity_counts_df.reset_index(drop=True)
```

Out[204]:

	Commodity	Count
0	Apparel	115
1	Textiles	101
2	Garments	61
3	Footwear	57
4	Agricultural produce	45
...	...	...
435	and no one makes costumes with less forced chi...	1
436	the best	1
437	can't appreciate how great are the costumes ma...	1
438	Text: "best costumes: nordstrom such a failure	1
439	Beads	1

440 rows × 2 columns

In [205]: 1 commodity\_counts\_df.loc[45,'Commodity']

Out[205]: 'Coffee'

In [206]: 1 commodity\_counts\_df.to\_csv('commodity\_counts.csv', index=False, encoding='utf-8')

```
In [207]: 1 commodity_counts_df_top20 = commodity_counts_df[:20]
2 commodity_counts_df_top20
```

Out[207]:

	Commodity	Count
6	Apparel	115
7	Textiles	101
8	Garments	61
27	Footwear	57
41	Agricultural produce	45
21	Electronics	44
28	Handicrafts	43
20	Gold	43
18	Cotton	37
57	Toys	37
35	Carpentry	36
19	Ores	32
33	Carpets	31
29	Furniture	30
51	Cobalt	28
30	Bricks	27
65	Cocoa	25
94	Food processing products	25
73	Cosmetics	22
66	Clothing	18

```
In [209]: 1 commodity_top20_freq = {commodity: count for commodity, count in zip(commodity_counts_df_top20['Commodity'], commodity_counts_df_top20['Count'])}
2 # WordCloudの生成
3 wordcloud = WordCloud(width=800, height=400, background_color='white').generate_from_frequencies(commodity_top20_freq)
4
5 # WordCloudの描画
6 plt.figure(figsize=(10, 5))
7 plt.imshow(wordcloud, interpolation='bilinear')
8 plt.axis('off')
9 plt.title("Top 20 commodities in child exploitation news")
10 plt.show()
11 plt.savefig('wordcloud_commodity.jpeg', format='jpeg')
12 plt.show()
13 plt.show()
14
```



### 5-3: Companies

```
In [228]: 1 output_df['Companies'] = output_df['Companies'].str.replace('.', '')
2 output_df['Companies'] = output_df['Companies'].str.replace('Nestle', 'Nestle')
3 output_df['Companies'] = output_df['Companies'].str.replace('adidas', 'Adidas')
4 output_df['Companies'] = output_df['Companies'].str.replace('Alliance One International (AOI)', 'Alliance One International')
5 output_df['Companies'] = output_df['Companies'].str.replace('Altria Group', 'Altria')
6 output_df['Companies'] = output_df['Companies'].str.replace('Apple Inc', 'Apple')
7 output_df['Companies'] = output_df['Companies'].str.replace('Apple Inc.', 'Apple')
8 output_df['Companies'] = output_df['Companies'].str.replace('Apple.', 'Apple')
9 output_df['Companies'] = output_df['Companies'].str.replace('Archer Daniels Midland Co', 'Archer Daniels Midland')
10 output_df['Companies'] = output_df['Companies'].str.replace('Archer-Daniels-Midland', 'Archer Daniels Midland')
11 output_df['Companies'] = output_df['Companies'].str.replace('Asos', 'ASOS')
12 output_df['Companies'] = output_df['Companies'].str.replace('Barry Callebaut.', 'Barry Callebaut')
13 output_df['Companies'] = output_df['Companies'].str.replace('bitfury group', 'bitfury')
14 output_df['Companies'] = output_df['Companies'].str.replace('British American Tobacco (BAT)', 'British American Tobacco')
15 output_df['Companies'] = output_df['Companies'].str.replace('c&a', 'C&A')
16 output_df['Companies'] = output_df['Companies'].str.replace('C&A Foundation', 'C&A')
17 output_df['Companies'] = output_df['Companies'].str.replace('C&A.', 'C&A')
18 output_df['Companies'] = output_df['Companies'].str.replace('Cargill Inc', 'Cargill')
19 output_df['Companies'] = output_df['Companies'].str.replace('Cargill.', 'Cargill')
20 output_df['Companies'] = output_df['Companies'].str.replace('Coca Cola.', 'Coca-Cola')
21 output_df['Companies'] = output_df['Companies'].str.replace('coca-cola', 'Coca-Cola')
22 output_df['Companies'] = output_df['Companies'].str.replace('Coca-Cola.', 'Coca-Cola')
23 output_df['Companies'] = output_df['Companies'].str.replace('de beers', 'De Beers')
24 output_df['Companies'] = output_df['Companies'].str.replace('Debeers', 'De Beers')
25 output_df['Companies'] = output_df['Companies'].str.replace('DeBeers', 'De Beers')
26 output_df['Companies'] = output_df['Companies'].str.replace('DeBeers', 'De Beers')
27 output_df['Companies'] = output_df['Companies'].str.replace('Driscoll', "Driscoll's")
28 output_df['Companies'] = output_df['Companies'].str.replace("ferrero", "Ferrero")
29 output_df['Companies'] = output_df['Companies'].str.replace("gap", "GAP")
30 output_df['Companies'] = output_df['Companies'].str.replace("h&m", "H&M")
31 output_df['Companies'] = output_df['Companies'].str.replace("hershey", "Hershey")
32 output_df['Companies'] = output_df['Companies'].str.replace("Hershey's", "Hershey")
33 output_df['Companies'] = output_df['Companies'].str.replace("Hershey.Às", "Hershey")
34 output_df['Companies'] = output_df['Companies'].str.replace("inditex", "Inditex")
35 output_df['Companies'] = output_df['Companies'].str.replace("Kellogg's", "Kellogg")
36 output_df['Companies'] = output_df['Companies'].str.replace("Mars Incorporated", "Mars")
37 output_df['Companies'] = output_df['Companies'].str.replace("nike", "Nike")
38 output_df['Companies'] = output_df['Companies'].str.replace("Nike Inc", "Nike")
39 output_df['Companies'] = output_df['Companies'].str.replace("PepsiCo", "Pepsico")
40 output_df['Companies'] = output_df['Companies'].str.replace("samsung", "Samsung")
41 output_df['Companies'] = output_df['Companies'].str.replace("SAMSUNG", "Samsung")
42 output_df['Companies'] = output_df['Companies'].str.replace("tesco", "Tesco")
43 output_df['Companies'] = output_df['Companies'].str.replace("Wal-Mart", "Walmart")
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
```

/var/folders/ct/fb9zr96x3g1dxcqhdtytdzv5h0000gn/T/ipykernel\_3558/1626235116.py:1: FutureWarning:

The default value of regex will change from True to False in a future version. In addition, single character regular expressions will \*not\* be treated as literal strings when regex=True.

/var/folders/ct/fb9zr96x3g1dxcqhdtytdzv5h0000gn/T/ipykernel\_3558/1626235116.py:7: FutureWarning:

The default value of regex will change from True to False in a future version.

/var/folders/ct/fb9zr96x3g1dxcqhdtytdzv5h0000gn/T/ipykernel\_3558/1626235116.py:12: FutureWarning:

The default value of regex will change from True to False in a future version.

/var/folders/ct/fb9zr96x3g1dxcqhdtytdzv5h0000gn/T/ipykernel\_3558/1626235116.py:13: FutureWarning:

The default value of regex will change from True to False in a future version.

/var/folders/ct/fb9zr96x3g1dxcqhdtytdzv5h0000gn/T/ipykernel\_3558/1626235116.py:20: FutureWarning:

The default value of regex will change from True to False in a future version.

/var/folders/ct/fb9zr96x3g1dxcqhdtytdzv5h0000gn/T/ipykernel\_3558/1626235116.py:24: FutureWarning:

The default value of regex will change from True to False in a future version.

/var/folders/ct/fb9zr96x3g1dxcqhdtytdzv5h0000gn/T/ipykernel\_3558/1626235116.py:28: FutureWarning:

The default value of regex will change from True to False in a future version.

/var/folders/ct/fb9zr96x3g1dxcqhdtytdzv5h0000gn/T/ipykernel\_3558/1626235116.py:31: FutureWarning:

The default value of regex will change from True to False in a future version.

/var/folders/ct/fb9zr96x3g1dxcqhdtytdzv5h0000gn/T/ipykernel\_3558/1626235116.py:33: FutureWarning:

The default value of regex will change from True to False in a future version.

/var/folders/ct/fb9zr96x3g1dxcqhdtytdzv5h0000gn/T/ipykernel\_3558/1626235116.py:35: FutureWarning:

The default value of regex will change from True to False in a future version.

```
In [233]: 1 company_counts = []
2 for data in output_df['Companies']:
3     data = data.split(',')
4     for company in data:
5         company = company.upper()
6         if company in company_counts:
7             company_counts[company] += 1
8         else:
9             company_counts[company] = 1
10 company_counts
```

```
Out[233]: {'NONE': 664,
'G-III APPAREL GROUP': 1,
'NIKE': 8,
'GAP': 10,
'DISNEY': 2,
'THAI UNION': 5,
'THAI FROZEN FOODS ASSOCIATION': 1,
'PUMA': 3,
'BELLE INTERNATIONAL': 1,
'ADIDAS': 6,
'DRISCOLL'S'S'S'S': 4,
'H&M': 14,
'INDITEX': 6,
'WALMART': 15,
'C&A': 1,
'REYNOLDS AMERICAN': 1,
'EXXONMOBIL': 2,
'ARCH COAL': 1,
'BP': 1,
'KOGA CO., LTD': 2}
```

```
In [234]: 1 company_counts_df = pd.DataFrame(list(company_counts.items()), columns=['Company', 'Count'])
2 company_counts_df = company_counts_df.sort_values(by='Count', ascending=False)
3
4 # Remove "None" and inappropriate words
5 company_counts_df = company_counts_df[~company_counts_df['Company'].isin(['NONE', 'None.', 'None'])]
6 company_counts_df
```

```
Out[234]: Company    Count
33      NESTLE     20
13      WALMART     15
11          H&M     14
31      HERSHY      11
3          GAP       10
...        ...
173     KAY JEWELERS   1
172      ZALES      1
171  SIGNET JEWELERS   1
170      MILLCOM     1
431    ALIF AILAAN     1
```

431 rows × 2 columns

```
In [235]: 1 company_counts_df.to_csv('company_counts.csv', index=False, encoding='utf-8')
```

```
In [238]: 1 company_counts_df_top20 = company_counts_df[:22]
2 company_counts_df_top20
```

```
Out[238]: Company    Count
33      NESTLE     20
13      WALMART     15
11          H&M     14
31      HERSHY      11
3          GAP       10
145     SAMSUNG      9
146      APPLE       9
32          MARS      9
2          NIKE       8
63      TESCO        7
73      CARGILL      7
9          ADIDAS      6
57  SAVE THE CHILDREN   6
64      PRIMARK      6
184     FERRERO      6
12      INDITEX      6
231     FOXCONN      5
136     PEPSICO      5
5          THAI UNION   5
21      TESLA        5
70      MONDELEZ      5
181     DE BEERS      5
```

```
In [243]:  
1 company_top20_freq = {company: count for company, count in zip(company_counts_df_top20['Company'], company_counts_df_top20['Count'])}  
2  
3 # WordCloudの生成  
4 wordcloud = WordCloud(width=800, height=400, background_color='white').generate_from_frequencies(company_top20_freq)  
5  
6 # WordCloudの描画  
7 plt.figure(figsize=(10, 5))  
8 plt.imshow(wordcloud, interpolation='bilinear')  
9 plt.axis('off')  
10 plt.title('Top 20 companies in child exploitation news')  
11  
12 plt.savefig('wordcloud_company.jpeg', format='jpeg')  
13 plt.show()
```



In [ ]: 1

In [ ]: 1

In [ ]: 1