# Report

**Web crawling task:**
- We had to crawl WikiCFP for big data, data mining, database,and AI conferences and their location every year.
- We wrote our own Java code and build the crawler based on that.
- When searching per category, WikiCFPallows navigation until page 20, so we crawled all 20 pages for big data, data mining,databases, and artificial intelligence. The output of the crawling was in the tab separated format: conference_acronym (event), conference_name(title), and conference_location(place).
- I had a separate output file for all the different categories.

**Data cleaning:**
- I used OpenRefine to clean each of these files (files for each category).
- I removed the unwanted fields, aggregated and segregated fields as per the requirement.
- I segregated the event acronym and the event year into two separate columns.
- I trimmed the country name from the location field as we only needed results according to the cities.
- Then I merged all the files together and created one single file of the crawled data.

## Splitting column1 to get the year column:



## Removing unwanted data:

Final data (after all the other operations):



**Tasks: In this part, we used Hadoop on the data we just crawled to compute various statistics.**

**Task I:** Compute and plot the number of conferences per city. Which are the top 10 locations?

- For this task, we wrote a Hadoop MapReduce program.
- Before that, we place the crawled data file in the HDFS using –put command.
- The driver class created a job and called the Mapper and Reducer classes.
- The Mapper class mapped all the similar cities together and passed it as the outputKey, along with an IntWritable variable (1) as outputValue which is used to sum up the count, to the reducer.
- The reducer then gets the city name and sum up the count 1 by 1 for each city.
- We export this package as a jar and use the following command to run it as Hadoop jar.

  - ```
    hadoop jar jarname.jar package.driver /inputpath /outputpath
    ```

Then we get the output file to our local file system or we can choose to cat the output from the Hadoop itself.

# Output: Top cities in descending order:



# Graph:

**Task II:** Output the list of conferences per city.

- For this task, I included created another Mapper and Reducer class.
- I created a separate job for these Mapper and Reducer class in the previous driver file itself.
- The Mapper passes each city as the outputKey to the reducer.
- Each conference then, is passes as the outputValue to the reducer, corresponding to its city.
- The reducer, in turn, appends each conference name to its corresponding city in which it happened.

**Output:** The list of city along with the conferences that took place in those cities.



**Task III:** For each conference regardless of the year (e.g., KDD), output the list of cities.

- I split the data that was the input for the mapper with regards to tabs.
- The conference name (event) was in the column 1 of my wifi_crawl data file.
- This I passed as the outputKey to the reducer.
- Then I passed the city name, which was column 4 of my data file, as outputValue to the reducer.
- The reducer then appended each of the conference name to its corresponding city.

**Output:** Conference with a list of cities in which they were held.



**Task IV:** For each city compute and plot a time series of number of conferences per year.

- For this task, I created 2 Mappers and 2 Reducers.
- The first mapper takes the event column, concatenates it with the year column value, and pass this new string as the outputValue.
- The city column is passed to the reducer as the ouputKey.
- The first reducer appends the list of events that happened in a year corresponding to the city.
- The second mapper then concatenates city with the year and passes it on as the outputKey.
- The outputValue is an IntWritable(1), which is used to increase the total count of each event that happened in that year in that particular city.

## Output:

Number of conferences per year that happened in a particular city.

```
Burlingame2012    1
Burlingame2013    1
Busan2012         1
CANCUN2017        1
Caen2017          1
Cagliari2017      1
Cairo2017         2
Cairo2018         1
Cala Millor2018 1
Calgary2018       1
Cali2018          1
Cambridge2012     1
Cambridge2017     3
Cambridge2018     1
Cancun2017        1
Casablanca2018    2
Central University of Technology2017      2
Certosa di Pontignano - Siena2018         1
Ceske Budejovice2018      1
Changa2017        2
Changchun2018     1
Changsha2014      1
Chania2012        3
ChengDu2017       2
Chengdu2016       4
Chengdu2018       18
Chennai2018       8
Chiang Mai2018 1
Chicago IL2014 1
Chicago2014       4
Chicago2017       1
Chongqing2018     4
Clermont-Ferrand2011      1
Coimbra2013       2
```

```
Jadavpur University2017 1
Jain College of Enginee2017       1
Jaipur2017        6
Jeju Island2017 1
Jeju Island2018 1
Jeju2018          1
Jinan2018         4
Jyvaskyla2018     1
KAIST2017         1
Kalamata2018      1
Kalbis Institute2017      1
Kanazawa2017      2
Kansas City2017 1
Kaohsiung2011     1
Kaohsiung2012     1
Kaohsiung2014     1
Kazan2013         1
Khenchela2016     1
Kitakyushu2017    1
Koblenz2013       1
Kohala Coast2015          1
Kolkata2017       1
Krakow2018        1
Kuala Lumpur2013          3
Kuala Lumpur2018          1
Kuantan2018       2
Kyoto Terrsa2013          1
Kyoto2017         2
Kyoto2018         4
La Habana2013     1
La Havana2018     2
La Rioja2017      1
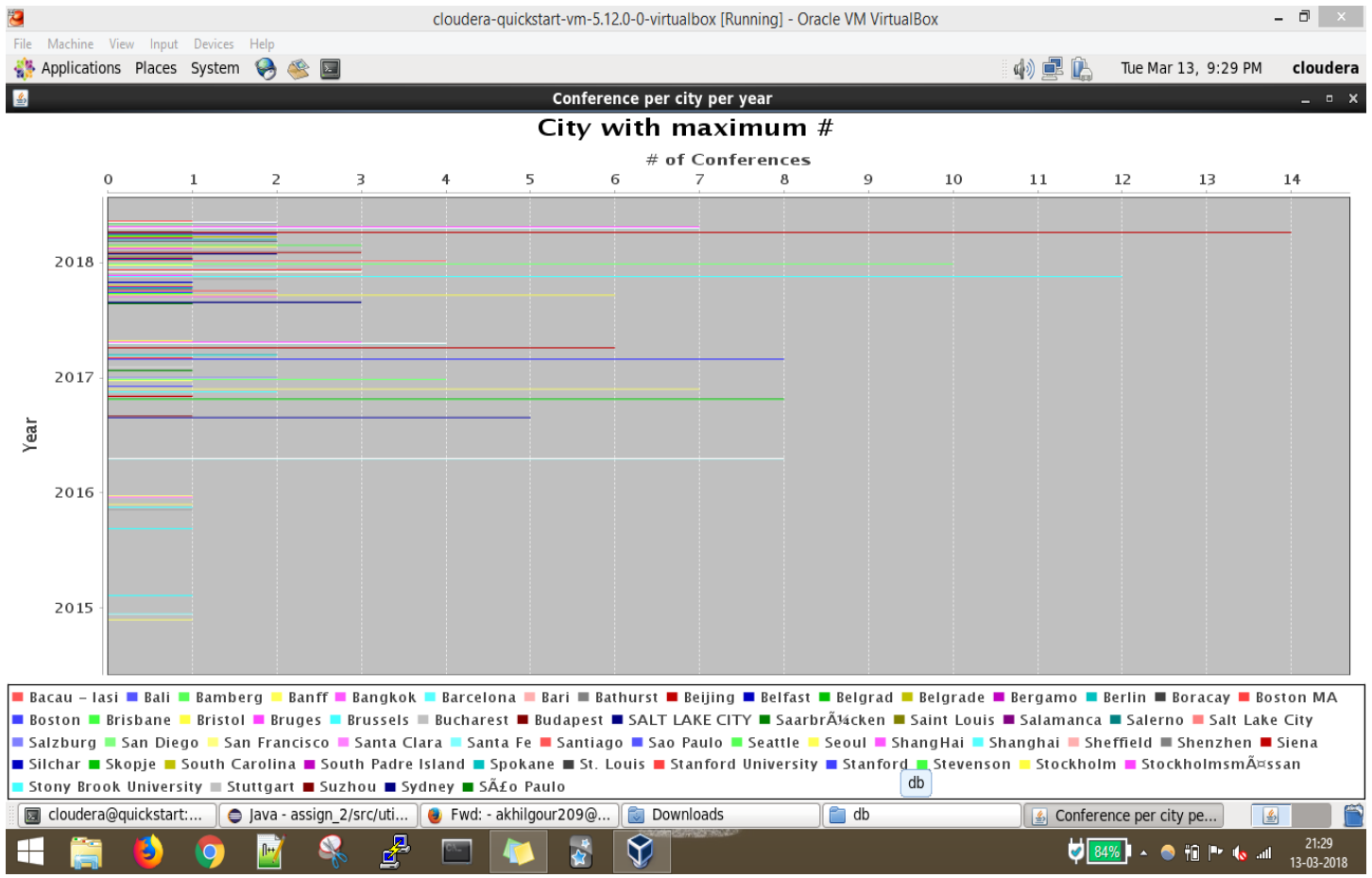Lago Di Garda2013         1
Laguna Hills2018          2
```

## Graph:



## Lessons learnt:

- How to use multiple Mapper/Reducer.
- Plotting instances of the output to a graph.
- How to include multiple Mapper/Reducer jobs within a java application.