

TASK 1

TASK 1.1: IMPLEMENT A PYTHON PROGRAM (EXTRACTION)

```
1  from pymongo import MongoClient
2
3  # Making database connection
4  client = MongoClient('localhost', 27017)
5  db = client["assignment1"]
6  col = db["movies"]
7  cur = col.find()
8
9  # Write a txt file with <year, company> pairs of all movies
10 file = open("year_and_company.txt", "w")
11 # Loop through each movie document in the cur
12 for movie in cur:
13     for i in range(min(len(movie['companies']), 3)):
14         year_company = movie['date'][-4:] + " , " + movie['companies'][i]['name'] + '\n'
15         file.write(year_company)
16 file.close()
```

Making database connection

Open a txt file for sorting year and company data

Loop through each movie document in the cur

For each I in the range from 0 to min(3, the length of movie['companies']):

 Extract the last four characters (year) of movie['date']

 Extract the company name from the 'companies' field

 Construct a string containing the year and company name using a delimiter “,”

 Split the line into year and company

Write the year_company file

Close the txt file

TASK 1.2: IMPLEMENT THE MAPREDUCE PROGRAM (COUNT)

```
1  from mrjob.job import MRJob
2
3  class MRWordCount(MRJob):
4      def mapper(self, _, line):
5          yield line, 1
6
7      def reducer(self, year_company, count):
8          yield year_company, sum(count)
9
10 if __name__ == "__main__":
11     MRWordCount.run()
```

Input the text file

Mapper function

Input: Each line contains a record in the format "year_company, count"

For each line in the input:

- Split the line into year_company and count using a delimiter ","

- Emit key-value pair (line, 1)

Reducer function

Input: List of key-value pairs

// key: year and company

//value: a list of counts

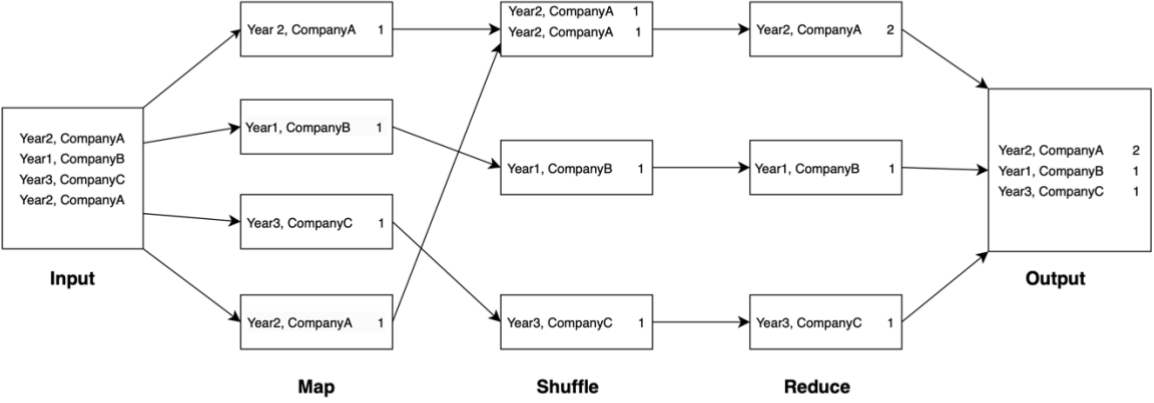
For each key-value pair in the input:

- Initialize a variable sum(count) to 0

- For each count in the list of counts:

- Emit key-value pair (year_company, sum(count))

FLOWCHART FOR TASK 1



TASK 2

TASK 2.1: MERGE SORT

```
1 from mrjob.job import MRJob
2 from mrjob.step import MRStep
3
4 class MRMergeSort(MRJob):
5
6     def steps(self):
7         return [
8             MRStep mapper=self.mapper, reducer=self.reducer
9         ]
10
11     def mapper(self, _, line):
12         year_company, count = line.split('\t')
13         count = int(count)
14         # "key" for combining year_company and count
15         yield "key", (year_company[1:-1], count)
16
17     def merge_sort(self, arr):
18         if len(arr) > 1:
19             mid = len(arr) // 2
20             left_half = arr[:mid]
21             right_half = arr[mid:]
22
23             self.merge_sort(left_half)
24             self.merge_sort(right_half)
25
26             i = j = k = 0
27             while i < len(left_half) and j < len(right_half):
28                 # Sort in ascending order
29                 if left_half[i][1] < right_half[j][1]:
30                     arr[k] = left_half[i]
31                     i += 1
32                 else:
33                     arr[k] = right_half[j]
34                     j += 1
35                 k += 1
36
37             while i < len(left_half):
38                 arr[k] = left_half[i]
39                 i += 1
40                 k += 1
41
42             while j < len(right_half):
43                 arr[k] = right_half[j]
44                 j += 1
45                 k += 1
46
47             # Return the sorted array
48             return arr
49
50     def reducer(self, _, tuples):
51         sorted_tuples = self.merge_sort(list(tuples))
52         for tuple in sorted_tuples:
53             yield tuple[0], tuple[1]
54
55 if __name__ == '__main__':
56     MRMergeSort.run()
```

Input the txt file

Steps function

Mapper function

Input: Each line contains a record in the format "year_company, count"

For each line in the input

 Split the line into year_company and count using a delimiter "\t"

 Combine year_company and count by using "key"

 Emit key-value pair ("key" (year_company, count))

yield "key", (year_company[1:-1], count)

Merge_sort function

If the length of array > 1:

 Calculate the middle index mid

 Split the array into left_half and right_half

 Call merge_sort on left_half and right_half

 Initialize variables i, j, and k to 0

 While i < length of left_half and j < length of right_half:

 Compare elements at indices i and j based on the second element of each tuple

 Update array at index k with element from left_half

 Increment i

 Else:

 Update array at index k with the element from right_half

 Increment j

 Increment k

 Copy any remaining elements from left_half

 Copy any remaining elements from right_half

Return the sorted array

Reducer function

Input: List of key-value pairs

// key: year_company

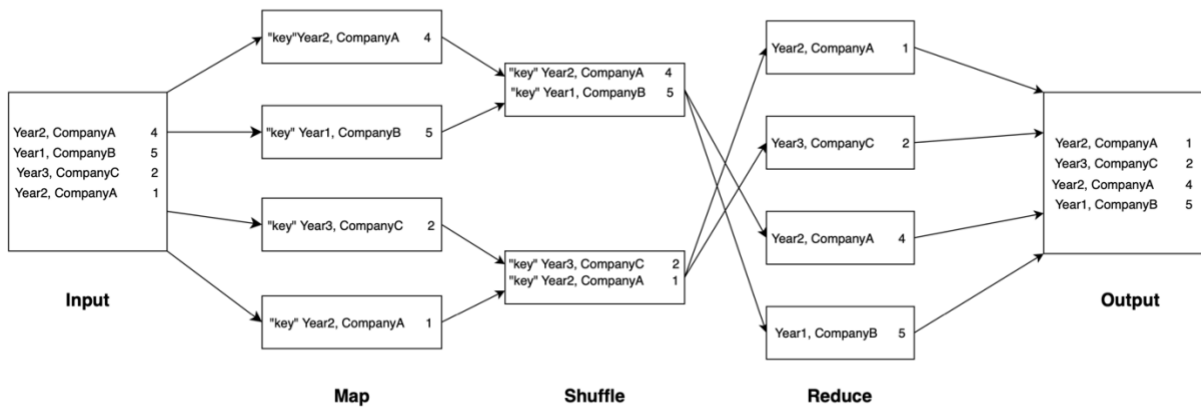
// value: count

Sort the list of tuples using merge-sort

For each tuple in the sorted list:

 Emit key-value pair with the first element of the tuple (year_company) and the second element (count)

FLOWCHART FOR MERGE SORT



TASK 2.2: BUCKET SORT

```
1  from mrjob.job import MRJob
2  from mrjob.step import MRStep
3
4  class MRBucketSort(MRJob):
5
6      def configure_args(self):
7          super(MRBucketSort, self).configure_args()
8          self.add_passthru_arg('--num_buckets', type=int, default=250)
9          self.add_passthru_arg('--bucket_size', type=int, default=3)
10
11     def steps(self):
12         return [
13             MRStep(
14                 mapper=self.bucket_assignment_mapper
15             ),
16             MRStep(
17                 reducer=self.bucket_sort_reducer
18             ),
19             MRStep(
20                 reducer=self.bucketid_sort_reducer
21             )
22         ]
23
24     def bucket_assignment_mapper(self, _, line):
25         year_company, count = line.split('\t')
26         count = int(count)
27         bucket_id = count // self.options.bucket_size
28         yield bucket_id, (year_company[1:-1], count)
29
30
31     def bucket_sort_reducer(self, bucket_id, records):
32         # Sort in descending order
33         sorted_records = sorted(records, key=lambda x: (-x[1], x[0]))
34         for record in sorted_records:
35             yield "key", (bucket_id, record)
36
37     def bucketid_sort_reducer(self, key, bucketid_records):
38         for value in sorted(bucketid_records, key=lambda x: x[0], reverse=True):
39             yield value[1]
40
41 if __name__ == '__main__':
42     MRBucketSort.run()
43
```

Input the txt file

Configure command-line arguments function

Steps function

Mapper function

Split the input line into year_company and count using a delimiter “\t”

Convert count to an integer

Calculate the bucket_id based on count and bucket size

Emit key-value pair

// key: bucket_id

// value: count

Reducer function for records (count) sorting

Sort records in descending order based on the second element of each tuple (records)

Emit sorted key-value pairs

Combine bucket_id and record by using “key”

// key: “key”

// tuple: bucket_id and record

Reducer function for bucket_id sorting

Sort records in descending order based on the first element of each tuple (bucket_id)

Emit sorted sorted values from bucketid_records

