Ministry of Higher Education and Scientific Research University Center of Mila Institute of Mathematics and Computer Science Department of Computer Science Master 2 I2A – Big Data 2025/2026

TD 2 – MapReduce Programming and Analysis

Objectives

Understand the MapReduce paradigm through practical data-processing exercises.

Exercise 1 – Word Frequency Counter

Goal: Count the occurrences of each word in a given text file.

Example Input:	Expected Output:
Hadoop is good for Big Data	Hadoop 1
Big Data means large data sets	is 1
	good 1
	for 1
	Big 2
	Data 2
	means 1
	large 1
	sets 1
Mapper logic:	
For each line in input:	
Split into words Emit (word, 1)	
Reducer logic:	
For each word:	
Sum all values Emit (word, total_count)	

Exercise 2 – Comparing Workload Growth

Goal: Compute the average temperature per weather station.

Example Input:	Expected Output:
ST01,2025-01-01,10	ST01 11
ST01,2025-01-02,12	ST02 8
ST02,2025-01-01,8	
Mapper logic: Emit(station, temperature)	·
Reducer logic: Compute average(temperature) for each station Emit(station, avg_temp)	

Exercise 3– Total Sales Per Product

Goal: Count the total revenue for each product

Example Input:	Expected Output:	
2025-01-01,Milk,20,100	Milk 3000	
2025-01-02,Milk,10,100	Bread 1000	
2025-01-01,Bread,50,20		
<pre>Mapper logic: product = fields[1] revenue = quantity * price Emit(product, revenue)</pre>		
Reducer logic: Sum all revenues for each product Emit(product, total_revenue)		

Exercise 4 – Average Movie Rating

Goal: Compute the average rating per movie

Example Input:	Expected Output:
M001,U01,4 M001,U02,5 M002,U01,3	M001 4.5 M002 3
Mapper logic: Emit(movie_id, rating)	
Reducer logic: Compute average rating per movie Emit(movie_id, avg_rating)	

Exercise 5 – Word Count by First Letter

Goal: Count how many words start with each alphabet letter

Example Input:	Expected Output:
Hadoop is Highly Helpful	Н 3
	I 1
	• •

Mapper logic: Emit(first_letter, 1)
Reducer logic: Sum counts by first_letter Emit(letter, total_count)