

# DISTRIBUTED SYSTEMS ASSIGNMENT REPORT

---



**Assignment ID:** 3

**Student Name:** 徐春晖 XU Chunhui

**Student ID:** 12110304

## DESIGN

---

### Structure

```
1 |code
2 |   parking_analysis.ipynb # main analysis notebook
3 |   plot.py # plot functions
4 |   requirements.txt
5 |
6 |data
7 |   .gitkeep
8 |   parking_data_sz.csv # origin data (may absent in submission)
9 |   parking_utilization.png # r5 subtask picture
10 |   r1.csv
11 |   r2.csv
12 |   r3.csv
13 |   r4.csv
14 |   r5.csv
```

I will use code and necessary comments to describe this part

# Data Clean

```
1     def __init__(self):
2         self.sc: SparkContext
3         self.spark: SparkSession
4         self.df: DataFrame
5         self.__tasks: tuple[DataFrame, DataFrame,
6                               DataFrame, DataFrame, DataFrame]
7         self.__init_data()
8
9     def __init_data(self):
10        """Initialize the data and create the dataframe"""
11
12        # Create a Spark context
13        sc = SparkContext.getOrCreate()
14
15        # Create a Spark session
16        spark = SparkSession \
17            .builder \
18            .appName("Parking Data Analysis") \
19            .getOrCreate()
20
21        # Load data from csv file
22        df = spark.read.csv(
23            path=ParkingDb.__get_data_file_path("parking_data_sz.csv"),
24            header=True,
25            schema=ParkingDb.schema
26        )
27
28        # filter out invalid data by checking out_time > in_time
29        cleaned_data = df.filter(
30            f.col("out_time") > f.col("in_time"))
31
32        # calculate parking time length and drop out_time column
33        df_convert_out = cleaned_data.withColumn(
34            "parking_time_length",
35            f.unix_timestamp("out_time") - f.unix_timestamp("in_time")
36        ).drop("out_time")
37
38        # cache the final dataframe
39        # df_convert_out = df_convert_out.cache()
40        df_final = df_convert_out
41
42        # assign to instance variables
43        self.sc = sc
44        self.spark = spark
45        self.df = df_final
46        self.__tasks = self.__execute()
```

# Task implementation logic

## Task 1



```
1 def __task1(self) -> DataFrame:
2     return self.df.groupBy("section") \
3         .agg(f.countDistinct("berthage").alias("count")) \
4         # .orderBy("section")
```

## Task 2



```
1 def __task2(self) -> DataFrame:
2     return self.df.select("berthage", "section").distinct() \
3         # .orderBy("section", "berthage")
```

## Task 3



```
1 def __task3(self) -> DataFrame:
2     return self.df.groupBy("section") \
3         .agg(f.avg("parking_time_length").cast("integer").alias("avg_parking_time")) \
4         # .orderBy("section")
```

## Task 4



```
1 def __task4(self) -> DataFrame:
2     return self.df.groupBy("berthage") \
3         .agg(f.avg("parking_time_length").cast("integer").alias("avg_parking_time")) \
4         .orderBy(f.col("avg_parking_time").desc())
```

## Task 5



```
1 def __task5(self) -> DataFrame:
2     # create a new dataframe with hourly time range
3     df_hours = self.df.withColumn(
4         "start_time",
5         f.date_trunc("hour", f.col("in_time"))
```

```

6         ).withColumn(
7             "end_time",
8             f.date_trunc("hour", f.col("in_time")) + f.expr("INTERVAL 1 HOUR")
9         )
10
11     # calculate total berthages for each section
12     total_berthages = self.df.select("section", "berthage").distinct() \
13         .groupBy("section").count().withColumnRenamed("count", "total_berthages")
14
15     # calculate hourly usage for each section
16     hourly_usage = df_hours.groupBy("start_time", "end_time", "section") \
17         .agg(f.countDistinct("berthage").alias("count"))
18
19     # join the two dataframes and calculate the percentage
20     result: DataFrame = hourly_usage.join(total_berthages, "section") \
21         .withColumn(
22             "percentage",
23             f.format_number(f.col("count") / f.col("total_berthages") * 100, 1)
24         ).select("start_time", "end_time", "section", "count", "percentage") \
25         # .orderBy("section", "start_time")
26
27     return result

```

## RUNNING RESULT

---

### Result Show

Could be found in notebook's output and `.csv` files.

## Spark job's DAG

### Task 1

# Details for Job 3

**Status:** SUCCEEDED

**Submitted:** 2024/12/12 19:39:34

**Duration:** 58 ms

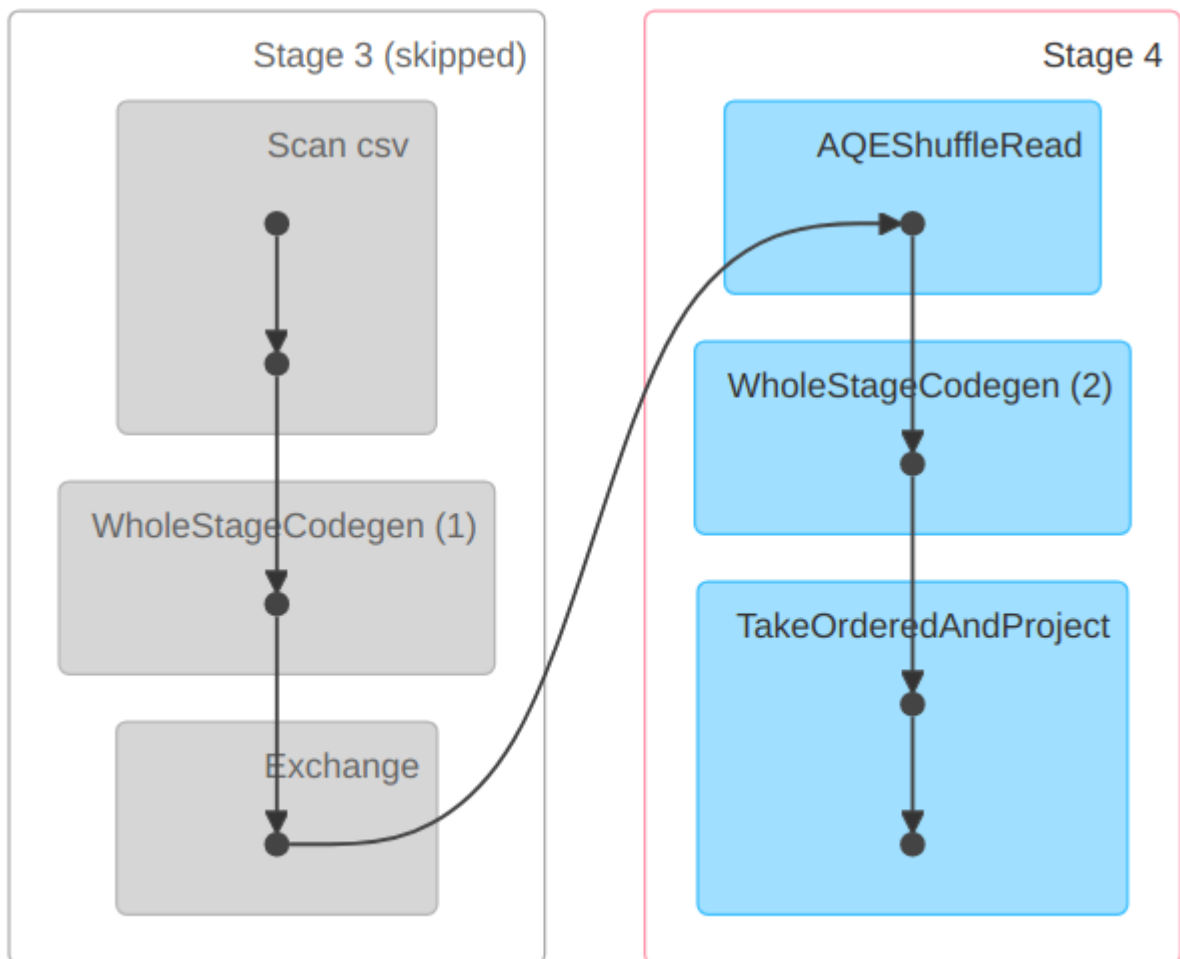
**Associated SQL Query:** 7

**Completed Stages:** 1

**Skipped Stages:** 1

► [Event Timeline](#)

▼ [DAG Visualization](#)



## Task 2

# Details for Job 7

**Status:** SUCCEEDED

**Submitted:** 2024/12/12 19:39:53

**Duration:** 38 ms

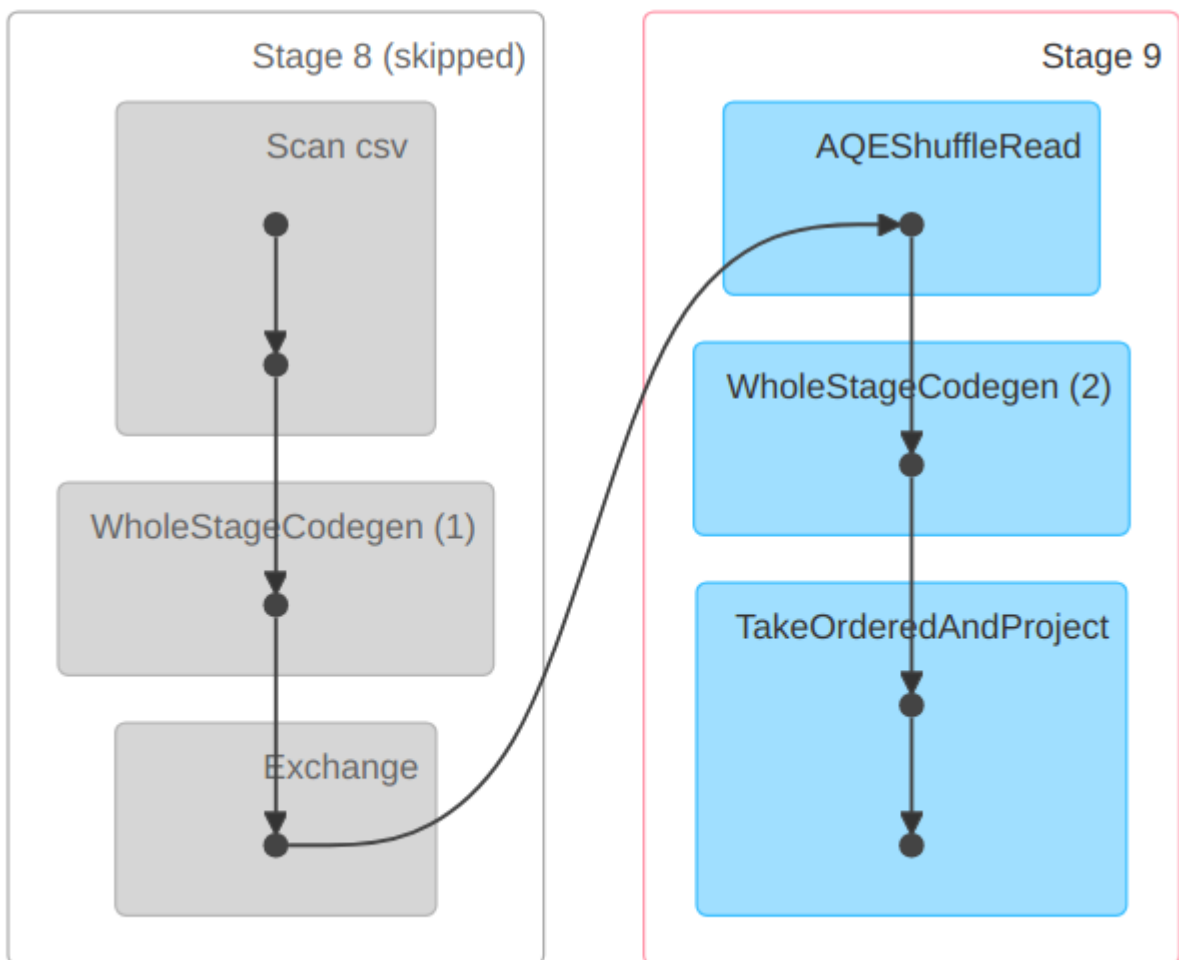
**Associated SQL Query:** 9

**Completed Stages:** 1

**Skipped Stages:** 1

► [Event Timeline](#)

▼ [DAG Visualization](#)

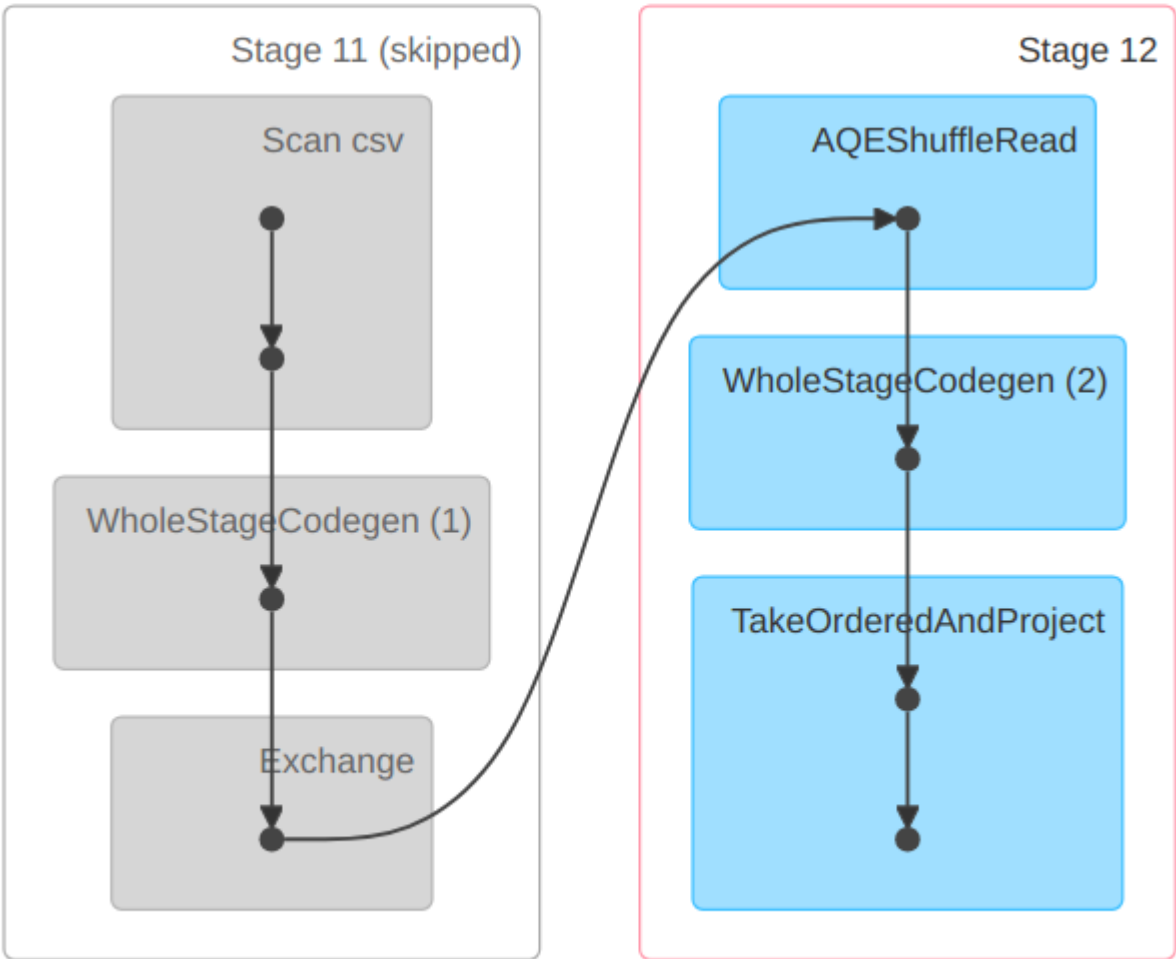


## Task 3

# Details for Job 9

Status: SUCCEEDED  
Submitted: 2024/12/12 19:39:54  
Duration: 87 ms  
Associated SQL Query: 10  
Completed Stages: 1  
Skipped Stages: 1

- ▶ Event Timeline
- ▼ DAG Visualization



## Task 4

# Details for Job 11

**Status:** SUCCEEDED

**Submitted:** 2024/12/12 19:39:55

**Duration:** 47 ms

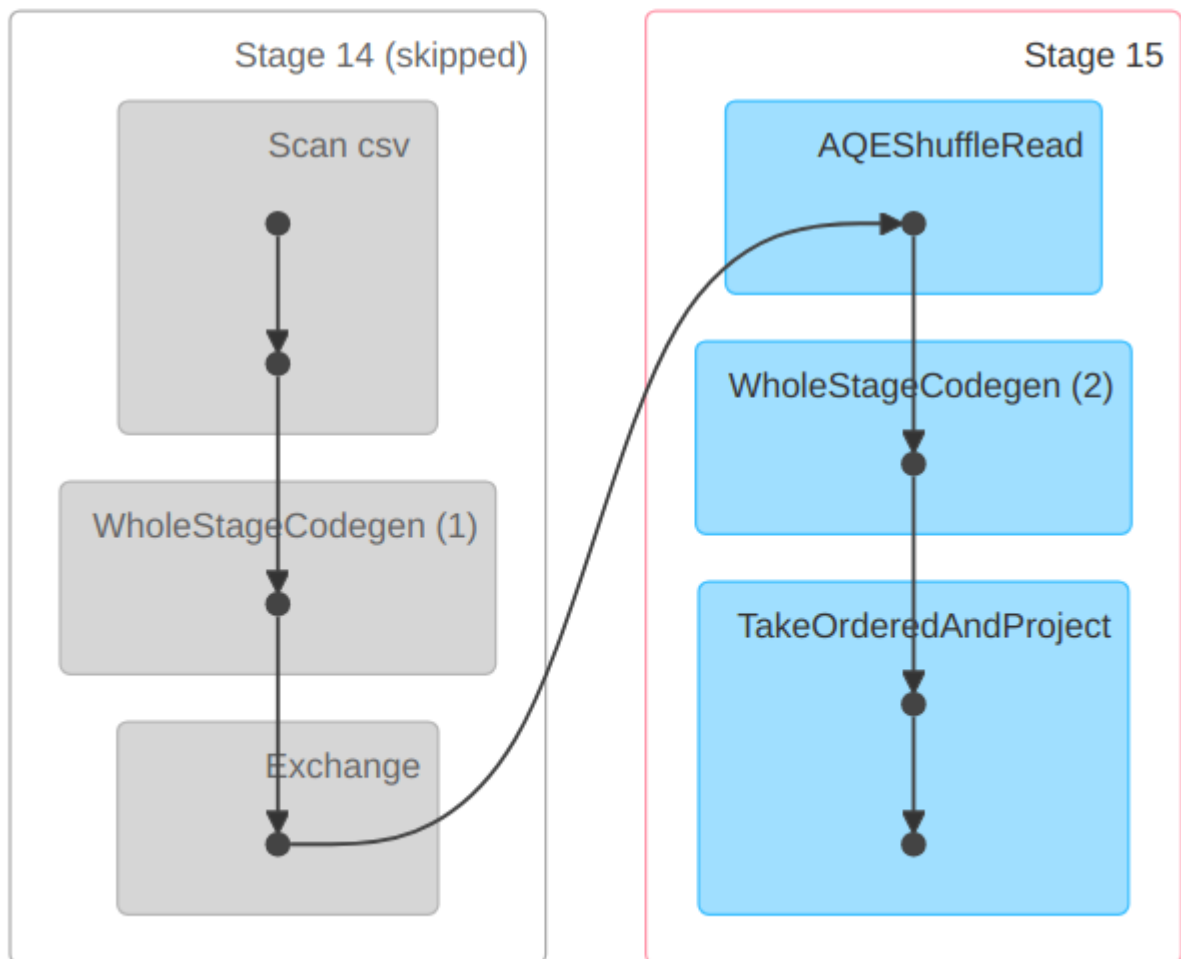
**Associated SQL Query:** [11](#)

**Completed Stages:** 1

**Skipped Stages:** 1

► [Event Timeline](#)

▼ [DAG Visualization](#)



## Task 5



# Details for Job 13

**Status:** SUCCEEDED

**Submitted:** 2024/12/12 19:39:57

**Duration:** 58 ms

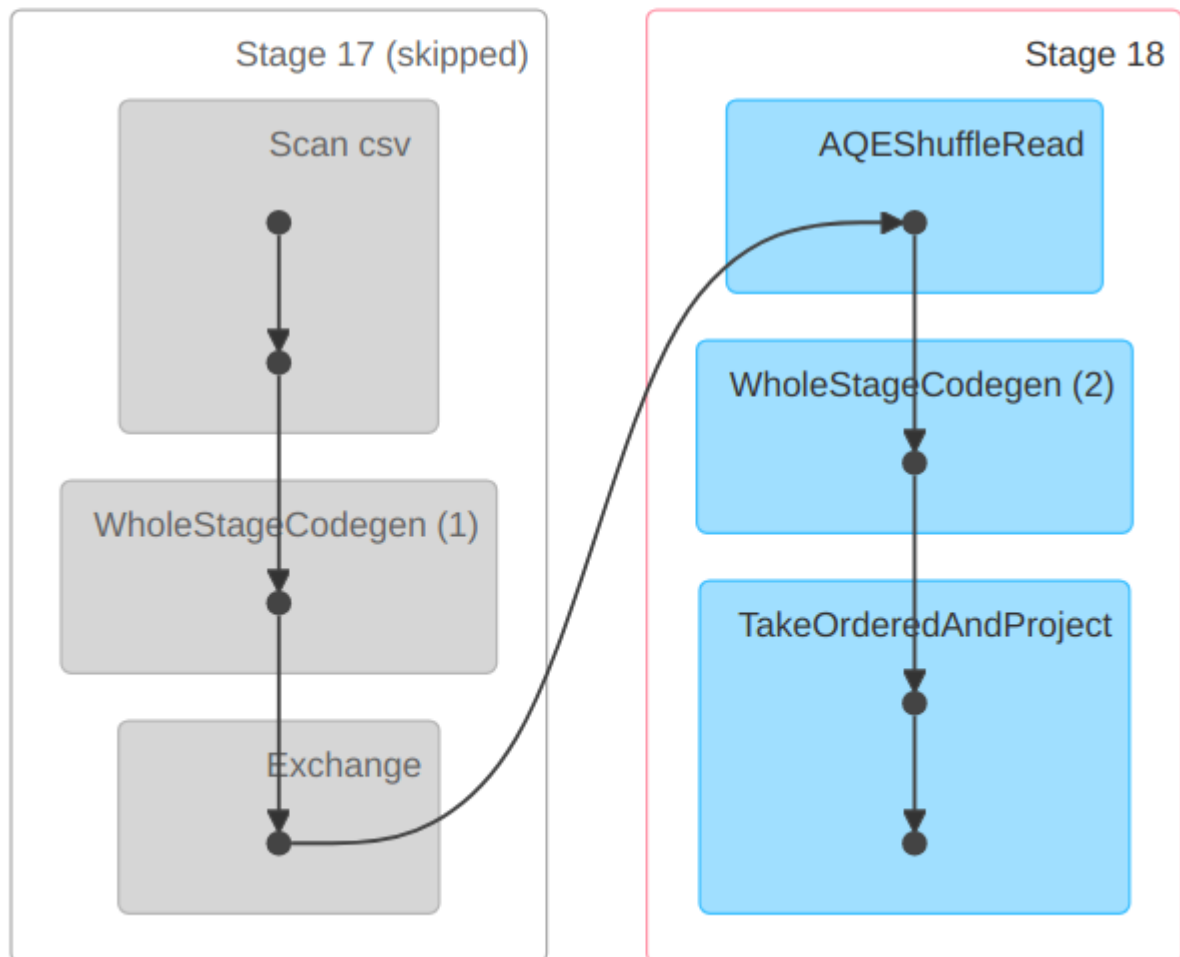
**Associated SQL Query:** [12](#)

**Completed Stages:** 1

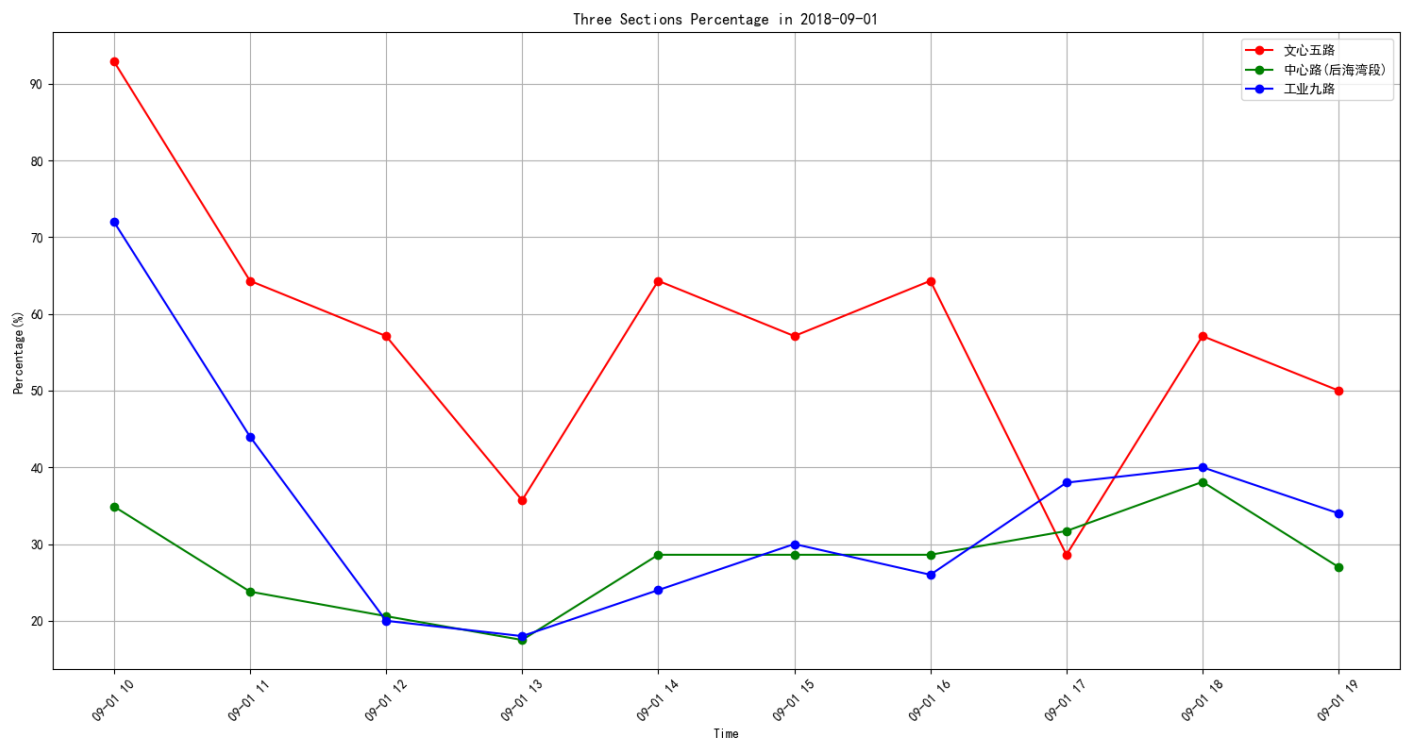
**Skipped Stages:** 1

► [Event Timeline](#)

▼ [DAG Visualization](#)



## r5 Subtask Analysis



I selected three sections and analyzed and visualized them hourly on 2018-09-01.

I found that parking occupancy rates were higher in the morning and evening, probably because people had not started their day yet, so the cars were in the parking lot.

During the day, people would drive around, so parking occupancy rates were lower.