

## Week 4

FIT5202 Big Data Processing

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



# Week 4 Agenda

- Week 3 Review
  - Review of Parallel Joins
  - Dataframe operations with pyspark
    - Sort
    - Distinct
    - Groupby
  - Dataframe UDFs(User Defined Functions)
  - TODO : Combining various operations to write dataframe queries.
  - Working on Assignment 1
- Assignment Project Exam Help
- <https://powcoder.com>
- Add WeChat powcoder

# Week 3 Review

- Spark Join Strategies
    - Broadcast Hash Join
    - Sort Merge Join
    - Shuffled Hash Join
  - Parallel Joins
    - Inner, Outer, Left, Right, Left Anti, Left Semi
  - Understanding Query Execution plans from Spark UI DAG
- Assignment Project Exam Help
- <https://powcoder.com>
- Add WeChat powcoder

# Week 3 Review

Why Full outer join uses “sort-merge”?

BHJ is not supported for full outer join.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Ref : <https://sujithjay.com/spark/broadcast-joins>

# PySpark Cheatsheet

[https://s3.amazonaws.com/assets.datacamp.com/blog\\_assets/PySpark\\_SQL\\_Cheat\\_Sheet\\_Python.pdf](https://s3.amazonaws.com/assets.datacamp.com/blog_assets/PySpark_SQL_Cheat_Sheet_Python.pdf)

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

# Spark Execution Plan

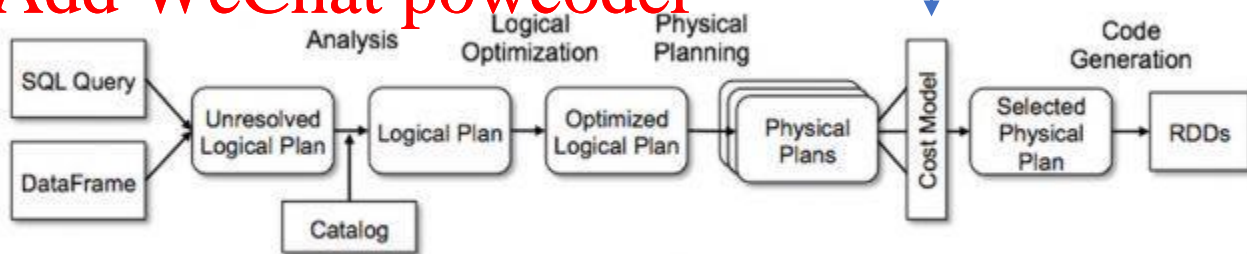
- Logical Plan
- Optimized Logical Plan
- Physical Plans
- Cost Model
- Selected Physical Plan

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Cost Based Optimization  
(time and resource taken by  
each strategy)



# Lab Instructions and Demo

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

# Thank You!

See you next week.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder