# Spark under the hood

02807 Computational Tools for Data Science

# Today

- The anatomy of a Spark job
- Advanced Spark functions
  - Window functions
  - Sketches in Spark
  - The physical plan
  - Collecting/writing data
  - Broadcast joins
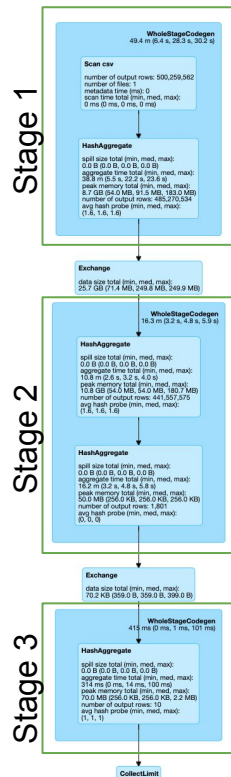- Project 3

# The anatomy of a Spark job

When you run an action in a cell in Jupyter:

- A **job** is submitted to Spark
- Data is split into logical **partitions**
- An execution **plan** is computed dividing your program into **stages**
- **Tasks** are executed on workers

A task is executing one stage on one partition of the data.

A job consists of more tasks than there are partitions, and there should be more partitions than cores in your system.

# The anatomy of a Spark job