

```
complaint_df = df.groupBy("Complaint Type").count().orderBy("count", ascending=False)

complaint_df.explain()

== Physical Plan ==
AdaptiveSparkPlan isFinalPlan=false
+- Sort [count#774L DESC NULLS LAST], true, 0
   +- Exchange rangepartitioning(count#774L DESC NULLS LAST, 200), ENSURE_REQUIREMENTS, [plan_id=586]
      +- HashAggregate(keys=[Complaint Type#608], functions=[count(1)])
         +- Exchange hashpartitioning(Complaint Type#608, 200), ENSURE_REQUIREMENTS, [plan_id=583]
            +- HashAggregate(keys=[Complaint Type#608], functions=[partial_count(1)])
               +- FileScan csv [Complaint Type#608] Batched: false, DataFilters: [], Format: CSV, Location: InMemoryFileIndex(1 paths)[file:/content/drive/MyDrive/Csv

complaint_df = df.groupBy("Borough").count().orderBy("count", ascending=False)

complaint_df.explain()

== Physical Plan ==
AdaptiveSparkPlan isFinalPlan=false
+- Sort [count#827L DESC NULLS LAST], true, 0
   +- Exchange rangepartitioning(count#827L DESC NULLS LAST, 200), ENSURE_REQUIREMENTS, [plan_id=606]
      +- HashAggregate(keys=[Borough#628], functions=[count(1)])
         +- Exchange hashpartitioning(Borough#628, 200), ENSURE_REQUIREMENTS, [plan_id=603]
            +- HashAggregate(keys=[Borough#628], functions=[partial_count(1)])
               +- FileScan csv [Borough#628] Batched: false, DataFilters: [], Format: CSV, Location: InMemoryFileIndex(1 paths)[file:/content/drive/MyDrive/Csv folder/
```

When working with big data, like the NYC 311 Service Requests dataset, it's easy to hit limits with traditional single-machine processing. These datasets can contain millions of rows and dozens of columns, and trying to load and process that much information on a single laptop can quickly lead to performance issues, or worse—crashes. This is where distributed processing becomes essential.

Distributed processing is basically the idea of splitting a large job across multiple computers (called nodes), which all work on smaller parts of the task at the same time. It's like splitting a big group project among friends instead of doing everything yourself—ideally, things get done faster and more efficiently. The key advantage here is parallelism. Instead of waiting for one machine to finish everything step-by-step, you get results from multiple workers running at once. That saves time and lets you work with much larger datasets than a single machine could handle.

Apache Spark is a great tool for distributed processing, especially in the context of data science. Spark keeps a lot of operations in memory. This means way less waiting around for data to load or write out between tasks, and in real-world terms, that means faster jobs. Even when I was working in Google Colab, I could still use Spark to load, clean, and process a large CSV in a way that would've been painful in regular Python using pandas.

Spark also gives you flexibility—it supports SQL-style queries, Data Frame APIs, and even lower-level RDDs for when you want more control. And because it's designed to scale, the same code can run on a laptop, a university cluster, or the cloud without major rewrites.

In short, distributed processing is crucial for handling big data, and Spark makes it more approachable and efficient. It's not perfect, and there's still a learning curve, but once you get it working, it feels like a major step up from the tools we used in earlier years.