# Persistence and Caching in Apache Spark

## Why Caching is Needed in Spark?

Caching (or persistence) in Apache Spark is an optimization technique that helps store intermediate results in memory, reducing re-computation and improving performance. It is especially useful in the following cases:

## 1. Repeated Use of Data

- When a dataset is used multiple times across transformations or actions, recomputing it each time can be expensive.
- Caching helps store the dataset in memory, so subsequent operations can directly access it without recalculating.

## 2. Reducing Disk Reads

- Reading data from disk is significantly slower than reading from memory.
- If data is stored on disk, each access involves expensive I/O operations, which can degrade performance.

## 3.Lazy Evaluation Impact

### 1. Repeated Computation in Lazy Evaluation

- Spark follows **lazy evaluation**, meaning it builds a **DAG (Directed Acyclic Graph)** of transformations but doesn't execute them until an **action** is triggered.
- When an action is finally called, Spark computes the entire lineage to produce the results.
- If multiple actions are executed on the same dataset, Spark will **re-run the entire computation each time**, leading to redundant processing.

### 2. Why Caching is Required?

- **Memory is not unlimited:**
  - Spark processes everything in memory, but if the dataset is too large to fit into memory, Spark may need to recompute parts of it and spill data to disk.
  - With caching, intermediate results are **stored efficiently**, reducing unnecessary re-computation.
- **Lazy evaluation and re-computation:**

- Without caching, Spark **recomputes everything** (entire lineage of transformations) whenever an action is triggered.
- With caching, the dataset is **stored after the first computation**, eliminating redundant recalculations.

**3.Reducing Disk & Network I/O Overhead**

Even though Spark operates in memory, large datasets may require:

1. **Reading from disk (HDFS, S3, etc.)** when memory is insufficient.
2. **Fetching data from remote nodes** in a distributed environment, which adds network latency.

# How and Where Caching Happens?

1. **Memory:**
   - Cached data is stored in a **deserialized format** for faster access.
   - If memory is insufficient, Spark **drops some cached data** to free up space.
2. **Disk:**
   - Data is **serialized and written (spill-over) to disk** when memory is full.
   - This helps when memory is limited but adds **disk I/O overhead**.

# When to Use Caching in Spark?

1. **When the dataset is used multiple times**
   - Without caching, Spark **recomputes** the dataset every time an action is triggered.
   - Caching **avoids redundant computation**.
2. **When data is not small enough to fit entirely in memory**
   - If memory is limited, caching helps manage intermediate results efficiently.
3. **Avoid caching excessively in memory-constrained environments**
   - Caching too much data in memory can degrade **overall cluster performance**.
4. **When performance is critical**
   - If reducing execution time is a priority, caching can significantly improve efficiency.
5. **When data is expensive to compute**
   - If transformations involve complex operations, caching saves processing time.
6. **When data does not change frequently**
   - If the dataset remains **static**, caching prevents unnecessary re-computation.