Let's clarify the **Disk I/O operations** for MapReduce in Hadoop, particularly in the case where **3 blocks** are involved. You're correct to seek more details on how to calculate read and write operations.

In a typical **MapReduce** job, there are several steps that result in **disk I/O** operations. For the **Map phase** and **Reduce phase**, let's break it down:

## Disk I/O in Hadoop (MapReduce)

1. **Map Phase:**

   - Each **Map task** processes a **block** of data (e.g., 128MB or 256MB block size).

   - The **Map task** reads the input data from the **DataNode**.

     - If there are **3 blocks** in total, there will be **3 read operations** — one for each block.

2. **Intermediate Write Operations:**

   - After reading and processing the block, the **Map task** generates **intermediate output** (key-value pairs).

   - These intermediate results are written to the **local disk** of the **Mapper node** (this is where disk I/O is happening).

   - If there are 3 map tasks (one per block), each will write its output, so there will be **3 write operations** for the intermediate results.

3. **Shuffle and Sort Phase:**

   - The **Shuffle and Sort** phase involves moving the intermediate results between nodes to group the data by keys for the **Reduce phase**. This is part of the MapReduce workflow but doesn't contribute directly to the read/write I/O from **DataNodes** because it's mostly about internal sorting and transferring data between the mappers and reducers.

4. **Reduce Phase:**

   - After the shuffle, **Reduce tasks** process the sorted data.

   - The **Reducer** writes the final output to the **HDFS**. If there are **3 reducers**, each will perform **write operations** for the final output.

# Corrected Disk I/O Calculation with Example

For **3 Blocks**, the following operations will occur:

- **Map Phase**:

    - **Reads**: 3 read operations (1 per block, since each block is processed by a Map task).

    - **Writes (Intermediate)**: 3 write operations (for the intermediate data produced by each Map task).

- **Reduce Phase**:

    - **Reads**: 3 read operations (each Reducer reads intermediate data).

    - **Writes (Final Output)**: 3 write operations (final output written to HDFS).

So, the total **Disk I/O** can be summarized as:

- **Reads**: 3 (for each block of data read during the Map phase) + 3 (for the intermediate data read by Reducers) = **6 reads**

- **Writes**: 3 (for intermediate output written during the Map phase) + 3 (for final output written during the Reduce phase) = **6 writes**

Thus, the total **Disk I/O** = **6 read operations** + **6 write operations** = **12 operations**.