

## What is a Memory Leak in Java (in Data Engineering)?

A **memory leak** in **Java** occurs when your **application keeps holding the memory (RAM)** that is **no longer needed**, and **does not release it** back to the system (JVM). This causes your application to **consume more and more memory** until it eventually **crashes** or slows down the performance.

In **Data Engineering**, where you are working with **large datasets, Spark, Kafka, Hadoop, or Big Data pipelines**, **memory leaks** can become a serious problem because:

- You often process **large volumes of data**.
- Heavy processing tasks require **memory-efficient operations**.
- If a memory leak happens, your application will **consume all available memory**, causing **OutOfMemoryError** or system crash.

## Why Do Memory Leaks Happen in Data Engineering (Java)?

### 1. Improper Use of Collections (HashMap, List, etc.)

- Suppose you have a **large dataset** and you store it in a **HashMap or ArrayList** without removing elements that are no longer needed.
- The **Java Garbage Collector (GC)** will not clear those objects because your code still holds a reference to them.
- This causes **memory leak** since old, unused objects still occupy memory.

### 2. Static Variables Holding Large Data

- If you use a **static variable** to store large data (like a **DataFrame, ArrayList, or ResultSet**), the **Garbage Collector (GC)** **cannot clear it** until the application is shut down.
- In **Data Engineering**, you often cache large data in memory. If not handled properly, it can lead to a **memory leak**.

### 3. Unclosed Resources (File, Database, Kafka Connection)

- In **data engineering**, you work with **files, databases, Kafka, Spark, or APIs**.
- If you open a **file, database connection, or Kafka consumer** without **closing it properly**, the memory will **never be released**, causing a memory leak.

### 4. Caching Large Data in Spark (Broadcast Variable Issue)

- In **Apache Spark (Data Engineering)**, if you use **broadcast variables** to share large datasets across nodes but **forget to clear them**, they will **stay in memory** until the application stops.
- This causes **memory leaks**.

## 5. Incorrect Usage of Threads in Data Processing

- In **data engineering**, you often use **multithreading** (for parallel processing).
- If you **start a thread but forget to stop it**, the thread holds memory and keeps running in the background, causing a memory leak.