

# External Memory Algorithms Report for DVD Rental Database Project

---

## 1. Introduction

In the context of database systems, external memory algorithms are essential when working with large datasets that cannot fit into the main memory. This is especially relevant in real-world applications like DVD rental databases, where large volumes of transactional data, film inventories, customer records, and rental histories are stored. Efficiently processing and sorting such large datasets is crucial for ensuring the system can handle data management tasks like sorting, searching, and updating records without significant delays.

In this report, we explore External Memory Algorithms in the context of a DVD Rental Database. The dataset in question includes tables for films, customers, rentals, payments, staff, and more, which are too large to fit into main memory. We focus on algorithms such as External Merge Sort that minimize disk I/O operations while processing large datasets, ensuring that database queries and updates remain efficient even when the data exceeds available RAM.

## 2. Memory Units and Dataset Characteristics

For this project, we use the following memory units:

- 1 Byte = 8 bits
- 1 Kilobyte (KB) = 1000 Bytes
- 1 Megabyte (MB) = 1000 Kilobytes (KB)
- 1 Gigabyte (GB) = 1000 Megabytes (MB)

### Dataset Characteristics:

The DVD rental database consists of several tables with different types of records. Some key tables and their characteristics are:

- Customer Table: Stores information about each customer (e.g., customer\_id, first\_name, last\_name, address, etc.)
- Film Table: Stores details of films available for rent (e.g., film\_id, title, description, release\_year, etc.)
- Rental Table: Records each rental transaction (e.g., rental\_id, rental\_date, customer\_id, inventory\_id, return\_date, etc.)
- Payment Table: Contains information about payments made by customers (e.g., payment\_id, amount, rental\_id, payment\_date, etc.)

- Staff Table: Details of the staff handling rentals and payments (e.g., staff\_id, first\_name, last\_name, store\_id).

Estimated Size of the Dataset:

Total number of records across all tables: 3 million records

Average size of each record: 500 bytes

Total size of the database:  $3,000,000 \times 500 \text{ bytes} = 1.5 \text{ GB}$

### 3. External Memory Algorithms for DVD Rental Database

#### 3.1 External Merge Sort

External Merge Sort is an ideal algorithm for sorting large datasets that do not fit entirely into memory. It works by dividing the dataset into smaller, manageable blocks, sorting these blocks in memory, and then merging the sorted blocks iteratively until the entire dataset is sorted. This approach ensures that the algorithm only needs to read and write the data from/to the disk a minimal number of times.

Algorithm Steps:

1. Splitting Phase: Divide the dataset into smaller chunks (blocks) that fit into the available memory (RAM). Sort each block in memory and write it back to external storage.
2. Merging Phase: After the first phase, the sorted blocks are merged. During each pass, the algorithm reads multiple sorted blocks into memory, merges them, and writes the merged result back to external storage. This process is repeated until only one sorted block remains.

#### 3.2 Example: External Merge Sort for DVD Rental Database

For the DVD rental database, let's assume:

- Dataset Size: 1.5 GB
- Block Size: 100 MB (chosen as a reasonable transfer unit between memory and disk)
- Number of Blocks:  $(1.5 \text{ GB} / 100 \text{ MB/block}) = 15 \text{ blocks}$

With 15 blocks, the External Merge Sort algorithm proceeds as follows:

1. First Pass (Splitting): The algorithm reads all 15 blocks, sorts them in memory, and writes the sorted blocks back to external storage. This involves:
  - 15 blocks read from disk.
  - 15 blocks written to disk.

2. Subsequent Passes (Merging): To merge the sorted blocks, assume the system can hold 5 buffers in memory (i.e., it can merge up to 4 blocks at once). The number of passes needed for merging is:

$$\lceil \log_5(15) \rceil = \lceil 2.4 \rceil = 3 \text{ passes.}$$

In each pass, the algorithm reads 5 blocks, merges them, and writes the merged result back to disk.

3. Total I/O Operations: The total I/O operations consist of reading and writing blocks during both the splitting and merging phases:

- First Pass: 15 blocks read and 15 blocks written.
- Subsequent Passes: 3 passes of reading and writing 15 blocks.

Total I/O operations:  $15 \text{ blocks} \times 2 \text{ (read + write)} \times 3 \text{ passes} = 90 \text{ I/O operations}$ .

## 4. Challenges in External Memory Algorithms

Several challenges arise when using external memory algorithms for large datasets like those in the DVD rental project:

1. Disk I/O Bottleneck: Disk I/O operations are much slower than in-memory operations, so minimizing disk access is a key challenge.
2. Data Locality: Poor data locality, such as non-contiguous data, can lead to increased disk seek times and slower processing.
3. Buffer Management: The number of buffers available in memory impacts how efficiently data can be processed. Allocating and managing buffers effectively is crucial for performance.
4. Merge Complexity: The process of merging sorted blocks can become complex as the number of blocks increases, requiring efficient strategies to manage disk accesses.

## 5. Performance Optimization Techniques

To optimize external memory algorithms, several techniques can be employed:

1. Buffer Size Optimization: Increasing the number of buffers used for merging allows more blocks to be processed in each pass, reducing the number of passes required.
2. Parallel Processing: If multiple storage devices or distributed systems are available, parallelizing the external memory algorithms can speed up the process.
3. Data Preprocessing: Preprocessing the data to ensure better data locality or sorting the data in smaller chunks can reduce disk I/O.
4. Efficient Merge Strategies: Using more sophisticated merge algorithms (e.g., multi-way merge) can reduce the overhead in each pass.

## 6. Conclusion

External memory algorithms like External Merge Sort are essential for managing large datasets in database systems, especially when the data exceeds available main memory. For the DVD Rental Database project, we demonstrated how External Merge Sort can efficiently sort a large dataset by minimizing disk I/O operations. The example calculations showed

the significant reduction in processing time when optimizing memory buffers and disk access patterns. By utilizing external memory algorithms and implementing performance optimizations, the DVD rental database system can efficiently handle large volumes of data and perform complex queries, updates, and sorts even when the data exceeds the capacity of the main memory.