

Project Context

Live services generate logs at a very high rate; e.g., our service creates over 100,000 log lines a second. Usually, these logs are loaded inside databases to enable fast querying, but the cost of keeping all the logs becomes too high. For this reason, only the recent logs are kept in databases, and logs for longer periods are kept in file archives. For this problem, we should assume we store our data in multiple files. We close a file and start a new file when the file size reaches 16GB. Our file names are of the format LogFile-#####.log (e.g., LogFile000008.log, or LogFile-000139.log). Currently, we have over 10,000 log files with the last log file name LogFile-0018203.log and a total data size of 285TB.

Problem Statement

We usually use our log database to query our logs. But now and then, we may have to query older logs for customer support or debugging. In most of these cases, we know the time range for which we need to analyze the logs. We need a tool that could extract the log lines from a given time range and print it to the console in a time effective manner. The command line (CLI) for the desired program is as below LogExtractor.exe -f "From Time" -t "To Time" -i "Log file directory location" All the time formats will be "ISO 8601" format. The extraction process should complete in a few seconds, minimizing the engineer's wait time

SOLUTION APPROACH :

USE MAIN MEMORY TO STORE RECENT TRANSACTIONS AND SECONDARY MEMORY FOR STORING EARLIER TRANSACTION AND LOAD NECESSARY TRANSACTION WHEN NEEDED SAME AS PAGE FAULT ALGORITHM IN OPERATING SYSTEM .

COMPLEXITY ANALYSIS

METHOD 1 : Assign timestamp to each transaction and apply binary search as given log will be sorted according to timestamp initially .

Complexity : $O(n \log n)$

METHOD 2 : Use Hashing to store logs as timestamp(key) and search in the hash table .

Average Complexity : $O(n)$

Worst case : $O(n^2)$ // will not happen if the hash function is implemented properly .