# Week#6 RocksDB Introduction

Sangeun Chae

2018314760

## 1.  INTRODUCTION

RocksDB 는 오픈 소스로 공개한 고속의 쓰기와 읽기에 최적화된 Key-Value 저장소이다. 멀티코어 환경의 SSD 저장 장치 기반의 서버 환경에서 상당한 성능을 보장하는 것이 특징이다. RocksDB 는 Log-structure merge-tree 를 바탕으로 설계된 데이터베이스 엔진이며, LSM tree 구조는, write 를 할 때 append only 방식으로 저장을 하기 때문에, write 를 sequential 하게 처리하여 B+Tree 에 비해 성능이 향상된다. 이번 랩에서는, RocksDB 의 설치와, DB_bench 의 실행을 목적으로 수행될 예정이다.

## 2.  METHODS

Github repository 에서 RocksDB 오픈 소스 코드를 클론 한 후에, 본인의 실험환경에서 RocksDB 를 make 한 후에, DB_bench 를 실행시킨다. DB_bench 를 실행시킨 후, 결과본을 저장한다.

## 3.  Performance Evaluation

### 3.1  Experimental Setup

| Type | Specification |
|---|---|
| OS | Ubuntu 20.04.3 LTS |
| CPU | AMD Ryzen 7 5800X 8-Core Processor (VMware support 4 Core) |
| Memory | 4GB |
| Kernel | Linux ubuntu 5.11.0.34-generic |
| Disk | VMware Virtual 80GB |

**Table 1: System setup**

| Type | Configuration |
|---|---|
| Bench Type | "readrandomwriterandom" |
| Direct flush_compaction | True |
| Direct read | True |
| Duration | 600 |

**Table 2: Benchmark setup**

## 3.2  Experimental Results



**Figure 1: Experimental result [1]**



**Figure 2: Experimental result [2]**

**Figure 3: Experimental result [3]**



**Figure 4: Experimental result [4]**



**Figure 5: Experimental result [5]**



**Figure 6: Experimental result [6]**

## 4. Conclusion

이번 랩에서는, RocksDB 의 설치와 DB_bench 의 실행을 목적으로 두고 진행했기 때문에, 조건의 변화없이 진행했다. 따라서, benchmark 를 실행했을 때, memtable hit count 와 cache hit count, 그리고 throughput 등의 결과 이외에는 유의미한 분석을 할 수 없었다.

## 5. REFERENCES

[1]  https://github.com/meeeejin/SWE3033-
     F2021/tree/main/week-6