

Introduction to Database Systems

Homework 2: Extendible Hashing

1. Introduction

In the lecture, you have learned what an index is in the database, the purpose of building an index, as well as a quick search in the database, which is hashing. In this homework, you are required to implement the extendible hashing and use it as a simple index for given data.

To assist you in better understanding this assignment, the following materials may be useful:

[Extendible Hashing \(Dynamic approach to DBMS\) - GeeksforGeeks](#)

[Extendible Hashing-A Fast Access Method for Dynamic Files](#)

Please continue reading for details.

2. Tasks

2.1 Task Description

This homework is divided into three parts, requiring you to implement basic functions such as **inserting data, querying, and deleting** of extendible hashing, and checking the correctness of the program through test data.

1. The data.txt, which contains two fields, key and value, as the input file for the implementation of an extendible hash table in C++.
2. The key_query.txt, which contains a field key. Based on the key order in this file, generate the corresponding value and local depth and export it as key_query_ans1.txt. If the value corresponding to the key is not found in the hash table, display the value as -1. However, the local depth should still be displayed.
3. The remove_key_query.txt, which contains a field key, to delete the corresponding keys in the constructed extendible hash table. Next, create a new file key_query_ans2.txt by reusing key_query.txt as input to generate the value and local depth again. If the value corresponding to the key is not found in the hash table, display the value as -1. However, the local depth should still be displayed.

You need to write two files, “**hash.h**” and “**hash.cpp**”, which are the header file and the cpp code of your hash-related classes. Please do not use the function like “map” or “unordered_map” to maintain the index.

2.2 File Descriptions

The provided files are explained below:

- `makefile`
 - TA will use provided makefile to compile your code when grading
- `main.cpp`
 - contains the main function which will be used to test your homework
 - main function will read all files needed for you, as well as record time and memory pages used by each task
 - **DO NOT** modify this file
- `utils.cpp`
 - contains some utility functions for file reading, used time output
 - **DO NOT** modify this file
- `utils.h`
 - header file for `utils.cpp`
 - **DO NOT** modify this file
- `data.txt`
 - contains the data used in the homework
 - each line consists of a key and a value, both are integers and are separated with a “,”
 - the data is not sorted on key or value
 - will be loaded as two vectors of integer in the main function
 - if the key exist more than one, take the last recorded value as the value in hash table
 - note that the data TAs use to test your program will not be identical to the one provided

- `key_query.txt`
 - contains the queries to test your extendible hash
 - each line consists of a key, which is an integer
 - will be loaded as a vector of int in the main function
 - some of the keys may not exist in the data
 - note that the queries TAs use to test your program will not be identical to the one provided
- `remove_key_query.txt`
 - contains the queries to be removed in extendible hash
 - each line consists of a keys, which is a integer
 - will be loaded as a vector of int in the main function
 - some of the key may not exist in the data
 - note that the queries TAs use to test your program will not be identical to the one provided
- `key_query_ans1.txt`
 - answer to `key_query.txt` right after `hash_table` constructed
 - start with a single line containing the global depth of the constructed hash table
 - record the value and the local depth of the bucket where the key is, both of which are integers and are separated with a “,”

```
3
1,3
2,3
3,3
```

- `key_query_ans2.txt`
 - answer to `key_query.txt` after removing query
 - start with a single line containing the global depth of the `hash_table` after removing query

- record the value and the local depth of the bucket where the key is, both of which are integers and are separated with a “,”

```
2
-1,1
-1,2
-1,1
```

You need to write two files, the detail requirements are explained below:

- hash.cpp
 - the code for your index, this file must contains a class named “**hash_table**”
 - hash_table class is an implementation of extendible hash index as explained in the slide.
 - At least four function needed to be implemented (add more if you need)
 - constructor hash_table()
 - key_query
 - remove_query
 - clear
 - hash_table(table_size, bucket_size, num_rows, key, value) takes five inputs, which are an integer initializing the hash table size(fixed to 2, so the initial global_depth = 1), an integer constraint about the size of the buckets(fixed to 3), an integer indicating the number of data rows, a vector of integer represent keys and a vector of integer represent values. You need to construct your hash index in this function by inserting the key, value pairs into the hash table **one by one**. While inserting, the global depth and the local depth should be maintained correctly.
 - key_query(query_keys, file_name) takes two inputs, which are a vector of integers indicating the key used for query and the name of the output fileThe function should generate an output file that starts with the current global depth of the hash table. Each subsequent row in the file should consist of two integers. The first integer is the value corresponding to the keys in query_keys (output -1 if the key is not found).The second integer is the local depth of the

bucket associated with the hash index derived from the keys (output the local depth the hash index need to be even if the key is not found).

- `remove_query` takes one input, which is a vector of integers indicating the keys to be removed. In this function you need to remove the key, value pair in the specific bucket. While removing, the global depth and the local depth should be maintained correctly. Note that the minimum value for both global depth and local depth is 1. Therefore, you do not need to merge buckets if their local depth is not greater than 1.
- `clear()` takes no input, the function frees all the memory used by your extendible hash index in this function.

- `hash.h`
 - header file for `hash.cpp`

We also provide the toy datasets in the `toy_dataset` directory for you to check your code's correctness. However, this dataset is not the test data we used.

2.3 Report file

For this homework, you should also submit a report file answering the following questions:

1. Explain your method for inserting a key in the hash table. How do you solve the overflow problem?
2. Explain your method for removing keys from the hash table. How do you shrink the bucket?
3. There are other methods for managing key-value pairs, such as B+ trees. Please explain the basic idea of a B+ tree and describe the differences between B+ trees and extendible hashing. Which one do you think is better? (There is no correct answer; feel free to express your idea.)

Requirements

1. You should implement the extendible hash table in C++.

2. Ensure the extendible hash table supports fundamental operations such as insertion, querying, and deletion, while maintaining global depth and local depth.
3. Test the program's accuracy using test data.
4. Submit the program's source code and the PDF report file

=====

Please note that the data files TAs use to test your program will be different to the ones provided, even in the number of rows and queries. Remember to make your program work with these differences in data.

If you successfully run the code and do the tasks, your program will output 3 files, “key_query_out1.txt”, “key_query_out2.txt” (generated by functions you implemented) and “time_used.txt” (automatically generated, you do not need to implement this). Below figure is the example way to examine your correctness.

```
(base) adsl-wcpeng@adslwcpeng-System-Product-Name:~/桌面/for_students$ make && ./hw2
make: 「hw2」 已是最新。
Data file reading complete, 1000000 rows loaded.
Key query file reading complete, 999999 queries loaded.
Remove query file reading complete, 999999 queries loaded.
(base) adsl-wcpeng@adslwcpeng-System-Product-Name:~/桌面/for_students$ diff key_query_ans1.txt key_query_out1.txt
(base) adsl-wcpeng@adslwcpeng-System-Product-Name:~/桌面/for_students$ diff key_query_ans2.txt key_query_out2.txt
```

3. Grading

In this homework, the correctness of your code only makes part of your score, to encourage you to optimize the index (e.g. concurrency, shrink), part of your score will be evaluated by the time used for the subtasks. The details are explained in the table below:

| Description | Score(%) |
|---|---------------|
| Correctness of “key_query_out1.txt” (global depth + value + local depth) | 5% + 5% + 15% |
| Correctness of “key_query_out2.txt” (global depth + value + local depth) | 5% + 5% + 15% |
| Correctness of extendible hash implementation | 20% |
| Time used to build index | 10% |
| Time used to process remove query | 10% |
| Report file | 10% |

For the 20% score related to the time performance of your code (time used to build index and time used to process remove query), we will rank the time taken for each subtask among all

students. Submissions that receive a wrong answer or have no submission will not be included in this ranking. The score will be awarded based on the percentile rank (PR) value.

- $PR \geq 90$: you get 100% of the performance score
- $60 \leq PR < 90$: you get 75% of the performance score
- $PR < 60$: you get 50% of the performance score
- If your answer/implementation is wrong, you will not get any points in the performance score

If you do not pass any “key_query_out1.txt” test data, you will not be able to get the score in “Time used to build index”. In the same way, if you do not pass any “key_query_out2.txt” test data, you will not be able to get the score in “Time used to process remove query”.

With the exception of “correctness of extendible hash implementation”, the homework will be revised automatically with a script. We will use your “hash.h” and “hash.cpp” with other provided codes to compile the program, we also provide the makefile we will use for your reference. The C++ version used will be C++17. Please make sure your code works with provided files and given settings.

4. Discussion

TAs had opened a channel **HW2 討論區** on Teams of the course, you can post questions about the homework on the forum. TAs will answer questions as soon as possible.

Discussion rules:

1. Do not ask for the answer to the homework.
2. Check if someone has asked the question you have before asking.
3. We encourage you to answer other students' questions, but again, do not give the answer to the homework. Reply the messages to answer questions.
4. Since we have this discussion forum, do not send email to ask questions about the homework unless the questions are personal and you do not want to ask publicly.

5. Submission

1. The deadline of this homework is **11/25 (Mon.) 23:55:00**.
2. The submission requires two code files (hash.h and hash.cpp) and a PDF report file
3. You should put all the files into one folder, the folder should be named as “**HW2_XXXXXXX**” where XXXXXXX is your student ID.
ex: **HW2_123456**

Then compress your folder into one `zip` file. Submit it to New E3 System with the format **HW2_XXXXXXX.zip** where XXXXXXX is your student ID.
ex: **HW2_123456.zip**

We **only accept one zip file**, each wrong format or naming format causes -10 points to your score (after considering late submission penalty).

4. Late submission lead to score of $(\text{original score}) \times 0.85^{\text{days}}$, for example, if you submit your homework right after the deadline, you get $(\text{original score}) \times 0.85$ points.
5. **0 points will be given to Plagiarism**. If the codes are similar to other people and you can't explain your code properly, you will be identified as plagiarism. TAs will strictly examine your code. It is important that you must write your code by yourself.
6. If there is anything you are not sure about submission, ask in the discussion forum.