

数据库课程设计 1: 可扩展哈希

1. 实验平台

Linux 或者 Windows

2. 编程语言

仅能使用 C/C++，不允许使用 STL。

3. 实现算法

可扩展哈希，具体算法请见课件<lec7 Hash-Based Indexes>中的第 5 至 12 张。

算法实现分为两大部分，第一部分是建立索引，第二部分是查询。建立索引是将输入的每一条记录根据指定的键值放入合适的哈希桶内，当哈希桶已满时，需要进行分裂。查询是根据输入的键值返回具有相同键值的记录，返回的记录可能有不止一条。

4. 实现过程

- (1) 读入由 tpc-h 生成的 lineitem.tbl，以 **L_ORDERKEY** 属性作为键值将记录放入合适的哈希桶内。
- (2) 读入测试文件 testinput.in 内的数据，数据中包含多个需要查询的键值，具体格式请见[6. 数据输入输出说明]。
- (3) 将通过键值查询得到的所有记录都输出到 testoutput.out 文件中，具体格式请见[6. 数据输入输出说明]。

5. 实现细节

- (1) 只能使用 **P 个页**，每个页的大小为 **8K bytes**，一个哈希桶的大小和一个页的大小相同。考虑以下两种情况：

- a) **P=8**（也就是整个内存中仅能使用 8 个页，这 8 个页用于存放索引和哈希桶的数据）
- b) **P=128**（也就是整个内存中仅能使用 128 个页，这 128 个页用于存放索引和哈希桶的数据）

由于页的数量有所限制，lineitem.tbl 的数据和对其建立的索引不可能都放在内存里，所以频繁的文件读写(I/O)是不可避免的。请比较 **P=8** 和 **P=128** 两种情况下，I/O 的次数、目录的大小、哈希桶的数量、查询的速度（每秒执行查询的数量）、I/O 用时占查询总用时的比例等。

- (2) 对于哈希桶内的数据，采用<**键，数据记录**>的方式进行存储。详细内容请参考课件<lec3 File and Indexing>中的第 16、17 张。
- (3) 对于存储哈希桶数据的页面，使用**变长记录**的方式进行存储。关于变长记录，请参考课件<lec 3 Storing Data>中的第 31 张。
- (4) 使用**时钟页面算法**进行页面置换，请参考<lec 3 Storing Data>中的第 19 张。
- (5) 分别实现**从低位和从高位进行扩展的哈希**。详细内容请参考课件<lec 7 Hash-Based Indexes>中的第 10 张。比较这两种哈希方法，包括桶的分裂方式、桶分裂时桶内数据的分配方式、I/O 的次数、目录的大小、桶的数量、查询的速度（每秒执行查询的数量）、I/O 用时占查询总用时的比例等。
- (6) 将建立的哈希索引输出到新的文件中，命名为 hashindex.out，格式自定。
- (7) 对于每个查询结果，请按照**L_PARTKEY**属性的值进行排序。

6. 数据输入输出说明

将文件路径作为启动参数传入程序，例如执行 `ExtendibleHash.exe` 程序， 执行命令为“`ExtendibleHash.exe D:\database`”。该路径表示 `lineitem.tbl` 和 `testinput.in` 所在的文件路径，同时也表示 `testoutput.out` 和 `hashindex.out` 输出的文件路径。以下是对文件的格式进行说明，请严格按照文件格式进行输入输出。`hashindex.out` 格式自定，因此不作说明。

(1) lineitem.tbl

通过 `tpc-h` 的 `dbgen` 程序生成, `Scale Factor` 设置为1, 所以记录的数量一定是 6001215。选择 `L_ORDERKEY` 属性作为键值建立哈希索引。

(2) testinput.in

数据第一行是 `n`，表示查询数量。接下来的 `n` 行，每行分别是一个整数，表示要查询的键值。

(3) testoutput.out

对于每一个查询，输出所有满足要求的记录，并且按照 `L_PARTKEY` 属性的值进行排序。每行一条记录。对于每个查询的结果，都以-1 作为结束。

以下是样例。

lineitem.tbl

```
lineitem.tbl
1 115519017706111721168.2310.0410.02N1996-03-131996-02-121996-03-22DELIVER IN PERSON|TRUCK|egular courts above the|
2 167310173112136145983.1610.0910.06N1996-04-121996-02-281996-04-20TAKE BACK RETURN|MAIL|ly final dependencies: slyly hold |
3 16370013701318113309.6010.1010.02N1996-01-291996-03-051996-01-31TAKE BACK RETURN|REG AIR|iously. regular, express dep|
4 121321463314128128955.6410.0910.06N1996-04-211996-03-301996-05-16NONE|AIR|lites. fluffily even de|
5 124027153415124122824.4810.1010.04N1996-03-301996-03-141996-04-01NONE|FOB| pending foxes. slyly re|
6 115635163816132149620.1610.0710.02N1996-01-301996-02-071996-02-03DELIVER IN PERSON|MAIL|arefully slyly ex|
7 210617011911138144694.4610.0010.05N1997-01-281997-01-141997-02-02TAKE BACK RETURN|RAIL|ven requests. deposits breach a|
8 314297117981145154058.0510.0610.00R1994-02-021994-01-041994-02-23NONE|AIR|ongside of the furiously brave acco|
9 3190361654012149146796.4710.1010.00R1993-11-091993-12-201993-11-24TAKE BACK RETURN|RAIL| unusual accounts. eve|
10 31284491347413127139890.8810.0610.07A1994-01-161993-11-221994-01-23DELIVER IN PERSON|SHIP|nal foxes wake. |
11 31293001893141212618.7610.0110.06A1993-12-041994-01-071994-01-01NONE|TRUCK|ly. fluffily pending d|
12 3180905165015128132986.5210.0410.00R1993-12-141994-01-101994-01-01TAKE BACK RETURN|FOB|ages nag slyly pending|
13 31621431966216126128733.6410.1010.02A1993-10-291993-12-181993-11-04TAKE BACK RETURN|RAIL|ges sleep after the caref|
14 4188035155601130130690.9010.0310.08N1996-01-101995-12-141996-01-18DELIVER IN PERSON|REG AIR|- quickly regular packages sleep. idly|
15 51108570185711115123678.5510.0210.04R1994-10-311994-08-311994-11-20NONE|AIR|ts wake furiously |
16 51239271392812126150723.9210.0710.08R1994-10-161994-09-251994-10-19NONE|FOB|sts use slyly quickly special instruc|
17 513753113513150173426.5010.0810.03A1994-08-081994-10-131994-08-26DELIVER IN PERSON|AIR|eodolites. fluffily unusual|
18 6139636121501137161998.3110.0810.03A1992-04-271992-05-151992-05-02TAKE BACK RETURN|TRUCK|p furiously special foxes|
19 718205219607112113608.6010.0710.03N1996-05-071996-03-131996-06-03TAKE BACK RETURN|FOB|ss pinto beans wake against th|
20 714524317758121911594.1610.0810.08N1996-02-011996-03-021996-02-19TAKE BACK RETURN|SHIP|es. instructions|
21 71947801979913146181639.8810.1010.07N1996-01-151996-03-271996-02-03COLLECT COD|MAIL| unusual reques|
22 71630731307416128131009.9610.0310.04N1996-03-221996-04-061996-04-20NONE|FOB|. slyly special requests hagg|
23 71518941944016139173943.8210.0810.01N1996-02-111996-02-241996-02-18DELIVER IN PERSON|TRUCK|ns haggle carefully ironic deposits. bl|
24 7179251175916135143058.7510.0610.03N1996-01-161996-02-231996-01-22TAKE BACK RETURN|FOB|sole. excuses wake carefully alongside of |
25 711572381226917516476.1510.0410.02N1996-02-101996-03-261996-02-13NONE|FOB|lthely regula|
26 32182704177211128147227.6010.0510.08N1995-10-231995-08-271995-10-26TAKE BACK RETURN|TRUCK|sleep quickly. req|
27 3219792144112132146605.4410.0210.00N1995-08-141995-10-071995-08-27COLLECT COD|AIR|lthely regular deposits. fluffily |
28 3214416166661312210.3210.0910.02N1995-08-071995-10-071995-08-23DELIVER IN PERSON|AIR| express accounts wake according to the|
29 321274317744141416582.9610.0910.03N1995-08-041995-10-011995-09-03NONE|REG AIR|e slyly final pac|
30 3218581183201514179059.6410.0510.06N1995-08-281995-08-201995-09-14DELIVER IN PERSON|AIR|symptotes nag according to the ironic depo|
31 321161514117161619159.6610.0410.03N1995-07-211995-09-231995-07-25COLLECT COD|RAIL| gifts cajole carefully.1|
32 33161396188551131140217.2310.0910.04A1993-10-291993-12-191993-11-08COLLECT COD|TRUCK|ng to the furiously ironic package|
33 331605191553212132147944.3210.0210.05A1993-12-091994-01-041993-12-28COLLECT COD|MAIL|gular theodolites|
34 3313746919963131517532.3010.0510.03A1993-12-091993-12-251993-12-23TAKE BACK RETURN|AIR|. stealthily bold exc|
35 3313918139191414175928.3110.0910.00R1993-11-091994-01-241993-11-11TAKE BACK RETURN|MAIL|unusual packages doubt caref|
```

testinput.in

```
testinput.in
1 4
2 1
3 4
4 10
5 32
```

testoutput.out

```
testoutput.out
1 121321463314128128955.6410.0910.06N1996-04-211996-03-301996-05-16NONE|AIR|lites. fluffily even de|
2 115635163816132149620.1610.0710.02N1996-01-301996-02-071996-02-03DELIVER IN PERSON|MAIL|arefully slyly ex|
3 124027153415124122824.4810.1010.04N1996-03-301996-03-141996-04-01NONE|FOB| pending foxes. slyly re|
4 16370013701318113309.6010.1010.02N1996-01-291996-03-051996-01-31TAKE BACK RETURN|REG AIR|iously. regular, express dep|
5 167310173112136145983.1610.0910.06N1996-04-121996-02-281996-04-20TAKE BACK RETURN|MAIL|ly final dependencies: slyly hold |
6 115519017706111721168.2310.0410.02N1996-03-131996-02-121996-03-22DELIVER IN PERSON|TRUCK|egular courts above the|
7 -1
8 4188035155601130130690.9010.0310.08N1996-01-101995-12-141996-01-18DELIVER IN PERSON|REG AIR|- quickly regular packages sleep. idly|
9 -1
10 -1
11 321274317744141416582.9610.0910.03N1995-08-041995-10-011995-09-03NONE|REG AIR|e slyly final pac|
12 321161514117161619159.6610.0410.03N1995-07-211995-09-231995-07-25COLLECT COD|RAIL| gifts cajole carefully.1|
13 3214416166661312210.3210.0910.02N1995-08-071995-10-071995-08-23DELIVER IN PERSON|AIR| express accounts wake according to the|
14 32182704177211128147227.6010.0510.08N1995-10-231995-08-271995-10-26TAKE BACK RETURN|TRUCK|sleep quickly. req|
15 3218581183201514179059.6410.0510.06N1995-08-281995-08-201995-09-14DELIVER IN PERSON|AIR|symptotes nag according to the ironic depo|
16 3219792144112132146605.4410.0210.00N1995-08-141995-10-071995-08-27COLLECT COD|AIR|lthely regular deposits. fluffily |
17 -1
```

7. 提交说明

请将所有要提交的内容放入一个压缩文件内，命名格式为“组号_组长姓名_实现平台”，例如“01_张三_linux”。组号信息请向学委咨询。要提交的内容如下：

(1) README

写明所提交的内容和小组信息。小组信息包括组员学号、姓名、班级。

(2) 源代码

放入“src”文件夹内。

(3) 可执行程序

放入“bin”文件夹内。

这里需要提交四个可执行程序。由于程序有两个参数，一个是页的数量，一个是哈希的扩展方式，根据排列组合，有 4 种方案。根据方案中参数的不同分别命名为 least_128, most_128, least_8, most_8。并附上 README 说明每个程序完成的功能和如何运行程序。

(4) 实验报告

可用中文或英文写。实验报告必须为 pdf 格式，请以“组号_report”的方式命名，例如“01_report”。

实验报告必须包含以下内容：

- a) 实验环境、工具说明
- b) 实验思路
- c) 数据结构设计
- d) 实验过程设计
- e) 代码结构
- f) 实验结果、不同参数下性能的比较（[5. 实验细节](1)和(5)中提到的比较内容）
- g) 心得体会

以上部分内容可根据需要使用简单的伪代码、流程图等表示。但切勿大篇幅粘贴代码！

可参考以下实验报告：

<http://ss.sysu.edu.cn/~fjl/sun.doc>

<http://ss.sysu.edu.cn/~fjl/ruan.pdf>

8. 注意

- (1) 务必将文件路径作为启动参数传入程序。
- (2) 请不要改动任何输入文件，即 lineitem.tbl 和 testinput.in。
- (3) 请严格按照要求进行输入输出和提交课程设计。
- (4) 请注意实验细节中提到的所有要求，请按照要求实现算法。对于没有明确要求的地方，可自行决定，并在实验报告中说明。
- (5) 建议不要使用 fscanf 和 fprintf，因为读写文件速度不理想。按照页的大小读入一定字节的数据会有利于提高 I/O 速度。
- (6) 虽然没有提供查询测试数据，但可自行设计查询测试数据来比较查询速度等。
- (7) 组内所有成员课程设计部分的得分都相同。
- (8) 如果某同学期末考试卷面分数低于 40 分，那么就可能将其课程设计分数乘以卷面分数，然后除以 100，作为最终的课程设计分数。

课程设计提交截止时间：2012 年 5 月 13 日 11:00 PM