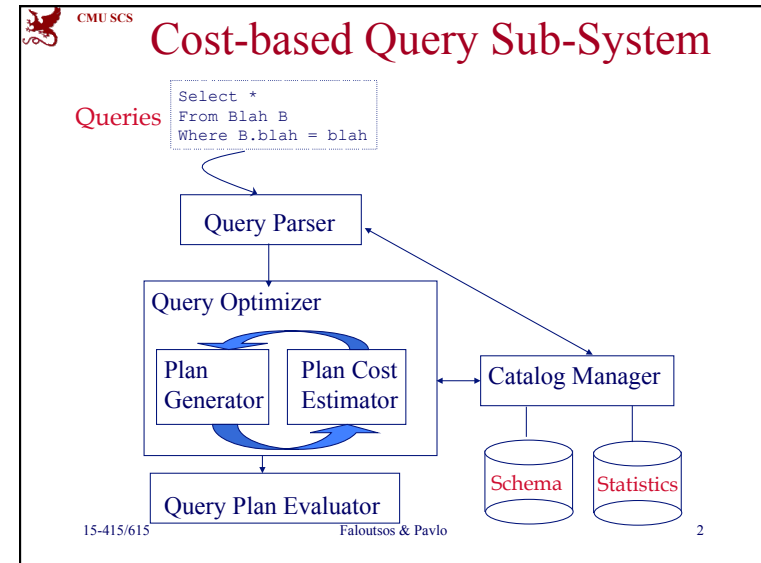


CMU SCS

Carnegie Mellon Univ.
Dept. of Computer Science
15-415/615 – DB Applications

Lecture #13: Query Evaluation
 (R&G ch. 12 and 14)

15-415/615 Faloutsos & Pavlo 1

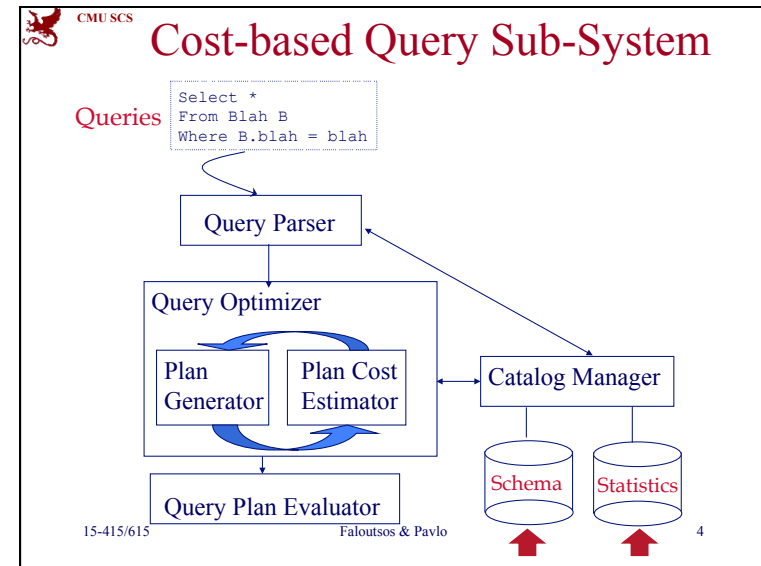



CMU SCS

Outline

- (12.1) Catalog
- (12.2) Intro to Operator Evaluation
- (12.3) Algo's for Relational Operations
- (12.6) Typical Q-optimizer
- (14.3.2) Hashing

15-415/615 Faloutsos & Pavlo 3






CMU SCS

Schema

- What would you store?

- How?

15-415/615 Faloutsos & Pavlo 5




CMU SCS

Schema

- What would you store?
- A: info about tables, attributes, indices, users
- How?
- A: in tables! eg.,
 - Attribute_Cat (attr_name: **string**, rel_name: **string**; type: **string**; position: **integer**)

15-415/615 Faloutsos & Pavlo 6




CMU SCS

Statistics

- Why do we need them?

- What would you store?

15-415/615 Faloutsos & Pavlo 7




CMU SCS

Statistics

- Why do we need them?
- A: To estimate cost of query plans
- What would you store?
 - NTuples(R): # records for table R
 - NPages(R): # pages for R
 - NKeys(I): # distinct key values for index I
 - INPages(I): # pages for index I
 - IHeight(I): # levels for I
 - ILow(I), IHigh(I): range of values for I

15-415/615 Faloutsos & Pavlo 8




CMU SCS

Outline

- (12.1) Catalog
- ➔ • (12.2) Intro to Operator Evaluation
- (12.3) Algo's for Relational Operations
- (12.6) Typical Q-optimizer
- (14.3.2) Hashing

15-415/615 Faloutsos & Pavlo 9




CMU SCS

Operator evaluation

3 methods we'll see often:

15-415/615 Faloutsos & Pavlo 10




CMU SCS

Operator evaluation

3 methods we'll see often:

- indexing
- iteration (= seq. scanning)
- partitioning (sorting and hashing)

15-415/615 Faloutsos & Pavlo 11




CMU SCS

``Access Path''

- Eg., index (tree, or hash), or scanning
- Selectivity of an access path:
 - % of pages we retrieve
- eg., selectivity of a hash index, on range query: 100% (no reduction!)

15-415/615 Faloutsos & Pavlo 12




CMU SCS

Outline

- (12.1) Catalog
- (12.2) Intro to Operator Evaluation
- ➔ • (12.3) Algo's for Relational Operations
- (12.6) Typical Q-optimizer
- (14.3.2) Hashing

15-415/615 Faloutsos & Pavlo 13




CMU SCS

Algorithms

- selection:
- projection
- join
- group by
- order by

15-415/615 Faloutsos & Pavlo 14




CMU SCS

Algorithms

- selection: scan; index
- projection (dup. elim.):
- join
- group by
- order by

15-415/615 Faloutsos & Pavlo 15



CMU SCS

Algorithms

- selection: scan; index
- projection (dup. elim.): hashing; sorting
- join
- group by
- order by

15-415/615 Faloutsos & Pavlo 16

CMU SCS

Algorithms

- selection: scan; index
- projection (dup. elim.): hashing; sorting
- join: many ways (loops, sort-merge, etc)
- group by
- order by

15-415/615 Faloutsos & Pavlo 17

CMU SCS

Algorithms

- selection: scan; index
- projection (dup. elim.): hashing; sorting
- join: many ways (loops, sort-merge, etc)
- group by: hashing; sorting
- order by: sorting

15-415/615 Faloutsos & Pavlo 18

CMU SCS

Iterator Interface

SELECT DISTINCT name, gpa
FROM Students

```

graph BT
    HeapScan((HeapScan)) -- "name, gpa" --> Sort((Sort))
    Sort -- "name, gpa" --> Distinct((Distinct))
  
```

15-415/615 Faloutsos & Pavlo 19

CMU SCS

Iterators

iterator

- Relational operators: subclasses of **iterator**:


```

class iterator {
    void init();
    tuple next();
    void close();
    iterator &inputs[];
    // additional state goes here
}
      
```
- iterators can be cascaded

15-415/615 Faloutsos & Pavlo 20

CMU SCS

Outline

- (12.1) Catalog
- (12.2) Intro to Operator Evaluation
- (12.3) Algo's for Relational Operations
- ➔ • (12.6) Typical Q-optimizer
- (14.3.2) Hashing

15-415/615 Faloutsos & Pavlo 21

CMU SCS

Q-opt steps

- bring query in internal form (eg., parse tree)
- ... into 'canonical form' (syntactic q-opt)
- generate alt. plans
- estimate cost; pick best

15-415/615 Faloutsos & Pavlo 22

CMU SCS

Q-opt - example

select name
from STUDENT, TAKES
where c-id= '415' and
STUDENT.ssn=TAKES.ssn

15-415/615 Faloutsos & Pavlo 23

CMU SCS

Q-opt - example

Canonical form

15-415/615 Faloutsos & Pavlo 24

CMU SCS

Q-opt - example

Hash join; merge join; nested loops;

Index; seq scan

STUDENT TAKES

15-415/615 Faloutsos & Pavlo 25

CMU SCS

Outline

- (12.1) Catalog
- (12.2) Intro to Operator Evaluation
- (12.3) Algo's for Relational Operations
- (12.6) Typical Q-optimizer
- ➔ • (14.3.2) Hashing

15-415/615 Faloutsos & Pavlo 26

CMU SCS

Grouping; Duplicate Elimination

select distinct ssn
from TAKES

- (Q1: what does it do, in English?)
- Q2: how to execute it?

15-415/615 Faloutsos & Pavlo 27

CMU SCS

An Alternative to Sorting: Hashing!

- Idea:
 - maybe we don't need the *order* of the sorted data
 - e.g.: forming groups in GROUP BY
 - e.g.: removing duplicates in DISTINCT
- Hashing does this!
 - And may be cheaper than sorting! (why?)
 - But what if table doesn't fit in memory??

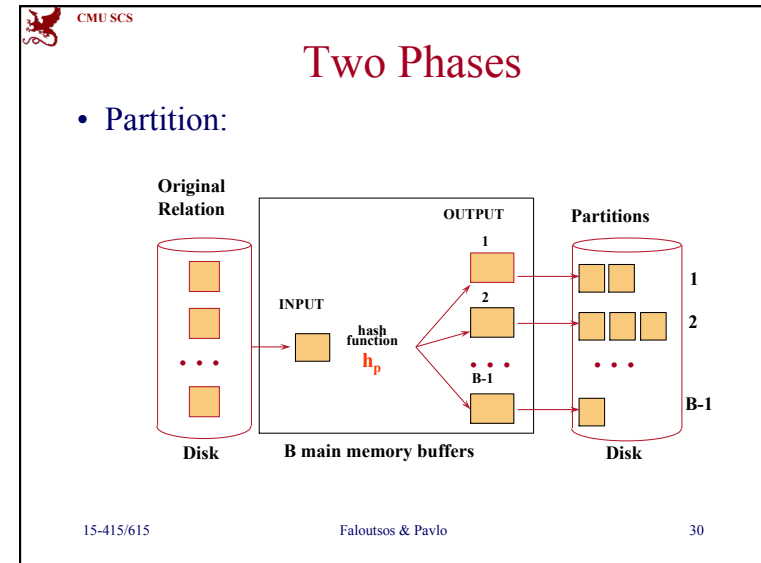
15-415/615 Faloutsos & Pavlo 28

CMU SCS

General Idea

- Two phases:
 - Phase 1: Partition:** use a hash function h_p to split tuples into partitions on disk.
 - We know that all matches live in the same partition.
 - Partitions are “spilled” to disk via output buffers

15-415/615 Faloutsos & Pavlo 29

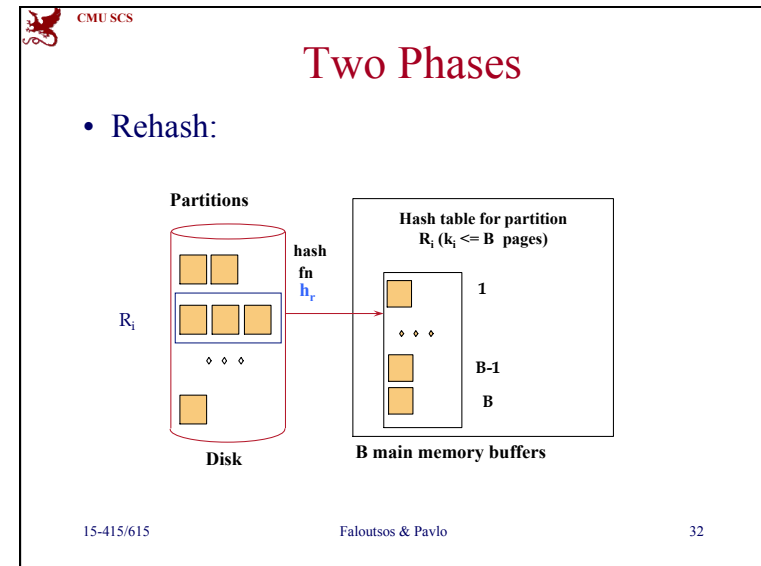



CMU SCS

General Idea

- Two phases:
 - Phase 2: ReHash:** for each partition on disk
 - (assuming it fits in memory)
 - read it into memory and build a main-memory hash table based on a hash function h_r
 - Then go through each bucket of this hash table to bring together matching tuples

15-415/615 Faloutsos & Pavlo 31






CMU SCS

Analysis

- How big of a table can we hash using this approach?
 - B-1 “spill partitions” in Phase 1
 - Each should be no more than B blocks big

15-415/615 Faloutsos & Pavlo 33




CMU SCS

Analysis

- How big of a table can we hash using this approach?
 - B-1 “spill partitions” in Phase 1
 - Each should be no more than B blocks big
 - Answer: $B(B-1)$.
 - ie., a table of N blocks needs about \sqrt{N} buffers
 - What assumption do we make?

15-415/615 Faloutsos & Pavlo 34




CMU SCS

Analysis

- How big of a table can we hash using this approach?
 - B-1 “spill partitions” in Phase 1
 - Each should be no more than B blocks big
 - Answer: $B(B-1)$.
 - ie., a table of N blocks needs about \sqrt{N} buffers
 - Note: assumes hash function distributes records **evenly!**
 - use a ‘fudge factor’ $f > 1$ for that: we need $B \sim \sqrt{f * N}$

15-415/615 Faloutsos & Pavlo 35

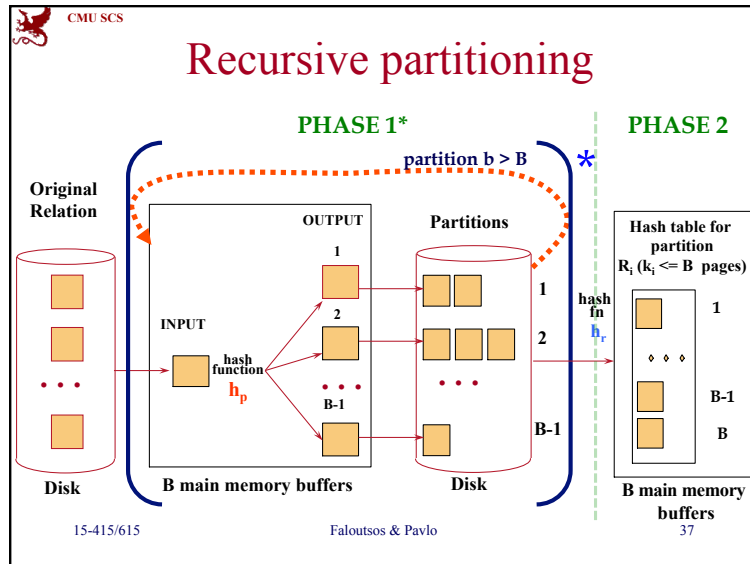


CMU SCS

Analysis

- Have a bigger table? **Recursive partitioning!**
 - In the ReHash phase, if a partition b is bigger than B, then recurse:
 - pretend that b is a table we need to hash, run the Partitioning phase on b , and then the ReHash phase on each of its (sub)partitions

15-415/615 Faloutsos & Pavlo 36



Real story break

- Partition + Rehash
- Performance is very slow!
- What could have gone wrong?

15-415/615 Faloutsos & Pavlo 38

Real story break

- Partition + Rehash
- Performance is very slow!
- What could have gone wrong?
- Hint: some buckets are empty; some others are way over-full.

15-415/615 Faloutsos & Pavlo 39

Hashing vs. Sorting

- Which one needs more buffers?

15-415/615 Faloutsos & Pavlo 40

CMU SCS

Hashing vs. Sorting

- **Recall: can hash a table of size N blocks in $\text{sqrt}(N)$ space**
- How big of a table can we sort in 2 passes?
 - Get N/B sorted runs after Pass 0
 - Can merge all runs in Pass 1 if $N/B \leq B-1$
 - Thus, we (roughly) require: $N \leq B^2$
 - We can sort a table of size N blocks in about space $\text{sqrt}(N)$
 - Same as hashing!

15-415/615 Faloutsos & Pavlo 41

CMU SCS

Hashing vs. Sorting

- Choice of sorting vs. hashing is subtle and depends on **optimizations** done in each case ...
 - Already discussed some optimizations for sorting:

15-415/615 Faloutsos & Pavlo 42

CMU SCS

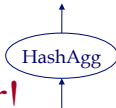
Hashing vs. Sorting

- Choice of sorting vs. hashing is subtle and depends on **optimizations** done in each case ...
 - Already discussed some optimizations for sorting:
 - (Heapsort in Pass 0 for longer runs)
 - **Chunk** I/O into large blocks to amortize seek+RD costs
 - **Double-buffering** to overlap CPU and I/O
 - Another optimization when using sorting for aggregation:
 - “Early aggregation” of records in sorted runs
 - We will discuss some optimizations for hashing next...

15-415/615 Faloutsos & Pavlo 43

CMU SCS

Hashing: We Can Do Better!



- Combine the summarization into the hashing process - How?

15-415/615 Faloutsos & Pavlo 44

CMU SCS

HashAgg

Hashing: We Can Do Better!

- Combine the summarization into the hashing process - How?
 - During the ReHash phase, don't store tuples, store pairs of the form **<GroupVals, RunningVals>**
 - When we want to insert a new tuple into the hash table
 - If we find a matching GroupVals, just update the RunningVals appropriately
 - Else insert a new <GroupVals, RunningVals> pair

```

select ssn, sum(credits)    (groupVal, runningVal)
from takes                 (12345, 12)
group by ssn               (45678, 18)
    
```

15-415/615 Faloutsos & Pavlo 45

CMU SCS

HashAgg

Hashing: We Can Do Better!

- Combine the summarization into the hashing process
- What's the benefit?
 - Q: How many entries will we have to handle?
 - A: Number of **distinct values** of GroupVals columns
 - Not** the number of tuples!!
 - Also probably "narrower" than the tuples

15-415/615 Faloutsos & Pavlo 46

CMU SCS

Even Better: Hybrid Hashing

- What if $B > \sqrt{N}$?
- e.g., $N=10,000$, $B=200$
- $B=100$ (actually, 101) would be enough for 2 passes
- How could we use the extra 100 buffers?

15-415/615

CMU SCS

Even Better: Hybrid Hashing

- What if $B > \sqrt{N}$?
- e.g., $N=10,000$, $B=200$
- $B=100$ (actually, 101) would be enough for 2 passes
- How could we use the extra 100 buffers?

A: 1ph for first partition;
2 for all others

15-415/615

CMU SCS

Even Better: Hybrid Hashing

- Idea: **hybrid!** ... keep 1st partition in memory during phase 1!
 - Output its stuff at the end of Phase 1.

15-415/615 Faloutsos & Pavlo 49

CMU SCS

Even Better: Hybrid Hashing

- What if $B=300$? (and $N=10,000$, again)
- i.e., 200 extra buffers?

15-415/615 Faloutsos & Pavlo 50

CMU SCS

Even Better: Hybrid Hashing

- What if $B=300$? (and $N=10,000$, again)
- i.e., 200 extra buffers?
- A: keep the first **2** partitions in main memory


15-415/615 Faloutsos & Pavlo 51

CMU SCS

Even Better: Hybrid Hashing

- What if $B=150$? (and $N=10,000$, again)
- i.e., 50 extra buffers?

15-415/615 Faloutsos & Pavlo 52




CMU SCS

Even Better: Hybrid Hashing

- What if $B=150$? (and $N=10,000$, again)
- i.e., 50 extra buffers?
- A: keep half of the first bucket in memory

15-415/615 Faloutsos & Pavlo 53




CMU SCS

Hybrid hashing

- can be used together with the summarization idea

15-415/615 Faloutsos & Pavlo 54




CMU SCS

So, hashing' s better ... right?

- Any caveats?

15-415/615 Faloutsos & Pavlo 55



CMU SCS

So, hashing' s better ... right?

- Any caveats?
- A1: sorting is better on non-uniform data
- A2: ... and when sorted output is required later.

Hashing vs. sorting:

- Commercial systems use either or both

15-415/615 Faloutsos & Pavlo 56



CMU SCS

Summary

- Query processing architecture:
 - Query optimizer translates SQL to a query plan
= graph of iterators
 - Query executor “interprets” the plan
- **Hashing** is a useful alternative to **sorting** for
dup. elim / group-by
 - Both are valuable techniques for a DBMS