

CMU SCS

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

C. Faloutsos – A. Pavlo
Lecture#13: Query Evaluation

CMU SCS

Administrivia

- HW4 is due this Thursday.
- Mid-term on **Tues March 3rd**
 - Will cover everything up to last week.
 - Closed book, one sheet of notes (double-sided)
 - Please email Christos + Andy if you need special accommodations.
 - More Info: <http://cmudb.io/s15-midterm>

Faloutsos/Pavlo CMU SCS 15-415/615 2

CMU SCS

Extended Office Hours

- **Christos:**
 - Wednesday Feb 25th 1:00pm-2:00pm
 - Friday Feb 27th 1:00pm-2:00pm
- **Andy:**
 - Wednesday Feb 25th 3:30pm-4:30pm
 - Monday Mar 2nd 1:00pm-2:00pm

Faloutsos/Pavlo CMU SCS 15-415/615 3

CMU SCS

Last Class

- **Sorting:**
 - External Merge Sort
- **Projection:**
 - External Merge Sort
 - Two-Phase Hashing

} **These are for when the data is larger than the amount of memory available.**

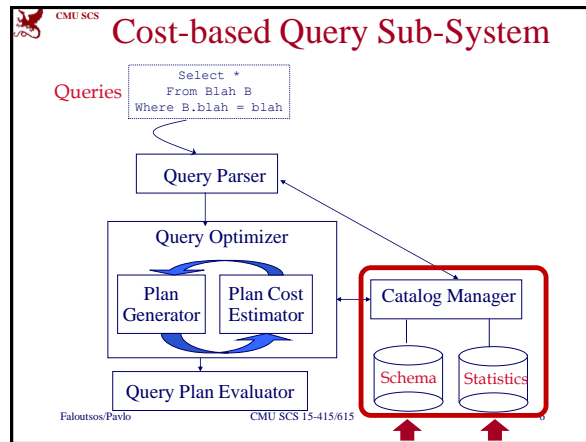
Faloutsos/Pavlo CMU SCS 15-415/615 4

CMU SCS

Today's Class

- Catalog (12.1)
- Intro to Operator Evaluation (12.2-3)
- Typical Query Optimizer (12.6)
- Projection/Aggregation: Sort vs. Hash (14.3.2)

Faloutsos/Pavlo CMU SCS 15-415/615 5



CMU SCS

Catalog: Schema

- What would you store?
- How?

Faloutsos/Pavlo CMU SCS 15-415/615 7

CMU SCS

Catalog: Schema

- What would you store?
 - Info about tables, attributes, indices, users
- How?
 - In tables!

```
Attribute_Cat (attr_name: string, rel_name:
string; type: string; position: integer)
```

Faloutsos/Pavlo CMU SCS 15-415/615 8

CMU SCS

Accessing Table Schema

- You can query the DBMS's internal **INFORMATION_SCHEMA** catalog to get info about the database.
- ANSI standard set of read-only views that provide info about all of the tables, views, columns, and procedures in a database
- Every DBMS also have non-standard shortcuts to do this.

Faloutsos/Pavlo CMU SCS 15-415/615 9

CMU SCS

Accessing Table Schema

- List all of the tables in the current database:

```
SELECT * FROM INFORMATION_SCHEMA.TABLES
WHERE table_catalog = '<db name>'
```

<code>\d;</code>	Postgres
<code>SHOW TABLES;</code>	MySQL
<code>.tables;</code>	SQLite

Faloutsos/Pavlo CMU SCS 15-415/615 10

CMU SCS

Accessing Table Schema

- List the column info for the student table:

```
SELECT * FROM INFORMATION_SCHEMA.COLUMNS
WHERE table_name = 'student'
```

<code>\d student;</code>	Postgres
<code>DESCRIBE student;</code>	MySQL
<code>.schema student;</code>	SQLite

Faloutsos/Pavlo CMU SCS 15-415/615 11

CMU SCS

Catalog: Statistics

- Why do we need them?
- What would you store?

Faloutsos/Pavlo CMU SCS 15-415/615 12

CMU SCS

Catalog: Statistics

- Why do we need them?
 - *To estimate cost of query plans*
- What would you store?
 - **Tables:** # tuples, # pages
 - **Indexes:** # distinct values, # pages, min/max

Faloutsos/Pavlo CMU SCS 15-415/615 13

CMU SCS

Today's Class

- Catalog (12.1)
- Intro to Operator Evaluation (12.2-3)
- Typical Query Optimizer (12.6)
- Projection/Aggregation: Sort vs. Hash (14.3.2)

Faloutsos/Pavlo CMU SCS 15-415/615 14

CMU SCS

Query Plan Example

```

SELECT cname, amt
FROM customer, account
WHERE customer.acctno =
      account.acctno
AND   account.amt > 1000
        
```

Relational Algebra:

$$\pi_{\text{cname, amt}}(\sigma_{\text{amt} > 1000}(\text{customer} \bowtie \text{account}))$$

Faloutsos/Pavlo CMU SCS 15-415/615 15

CMU SCS

Query Plan Example

```

SELECT cname, amt
FROM customer, account
WHERE customer.acctno =
      account.acctno
AND   account.amt > 1000
        
```

```

      graph BT
      C[CUSTOMER] --> J[⋈ acctno=acctno]
      A[ACCOUNT] --> J
      J --> S["σ amt > 1000"]
      S --> P["π cname, amt"]
        
```

Faloutsos/Pavlo CMU SCS 15-415/615 16

CMU SCS

Query Plan Example

```

SELECT cname, amt
FROM customer, account
WHERE customer.acctno =
      account.acctno
      AND account.amt > 1000
    
```

File Scan

File Scan

On-the-fly

On-the-fly

Nested Loop

π cname, amt

σ amt > 1000

acctno=acctno

CUSTOMER ACCOUNT

Faloutsos/Pavlo CMU SCS 15-415/615 17

CMU SCS

Query Plan Example

```

SELECT cname, amt
FROM customer, account
WHERE customer.acctno =
      account.acctno
      AND account.amt > 1000
    
```

The output of each operator is the input to the next operator.

Each operator iterates over its input and performs some task.

π cname, amt

σ amt > 1000

acctno=acctno

CUSTOMER ACCOUNT

Faloutsos/Pavlo CMU SCS 15-415/615 18

CMU SCS

Operator Evaluation

- Several algorithms are available for different relational operators.
- Each has its own performance trade-offs.
- The goal of the query optimizer is to choose the one that has the lowest “cost”.

Next Week: How the DBMS finds the best plan.

Faloutsos/Pavlo CMU SCS 15-415/615 19

CMU SCS

Operator Execution Strategies

- Indexing
- Iteration (= seq. scanning)
- Partitioning (sorting and hashing)

Faloutsos/Pavlo CMU SCS 15-415/615 20

CMU SCS

Operator Algorithms

- **Selection:**
- **Projection:**
- **Join:**
- **Group By:**
- **Order By:**

Faloutsos/Pavlo CMU SCS 15-415/615 21

CMU SCS

Operator Algorithms

- **Selection:** file scan; index scan
- **Projection:** hashing; sorting
- **Join:**
- **Group By:**
- **Order By:**

Faloutsos/Pavlo CMU SCS 15-415/615 22

CMU SCS

Operator Algorithms

- **Selection:** file scan; index scan
- **Projection:** hashing; sorting
- **Join:** many ways (loops, sort-merge, etc)
- **Group By:**
- **Order By:**

Faloutsos/Pavlo CMU SCS 15-415/615 23

CMU SCS

Operator Algorithms

- **Selection:** file scan; index scan
- **Projection:** hashing; sorting
- **Join:** many ways (loops, sort-merge, etc)
- **Group By:** hashing; sorting
- **Order By:** sorting

Faloutsos/Pavlo CMU SCS 15-415/615 24

CMU SCS

Operator Algorithms

- **Selection:** file scan; index scan
- **Projection:** hashing; sorting
- **Join:** many ways (loops, sort-merge, etc)
- **Group By:** sorting
- **Order By:** sorting

Faloutsos/Pavlo CMU SCS 15-415/615 25

CMU SCS

Today's Class

- Catalog (12.1)
- Intro to Operator Evaluation (12.2-3)
- Typical Query Optimizer (12.6)
- Projection/Aggregation: Sort vs. Hash (14.3.2)

Faloutsos/Pavlo CMU SCS 15-415/615 26

CMU SCS

Query Optimization

- Bring query in internal form (eg., parse tree)
- ... into "canonical form" (syntactic q-opt)
- ➔ Generate alternative plans.
- Estimate cost for each plan.
- Pick the best one.

Faloutsos/Pavlo CMU SCS 15-415/615 27

CMU SCS

Query Plan Example

```

SELECT cname, amt
FROM customer, account
WHERE customer.acctno =
      account.acctno
AND account.amt > 1000
    
```

Faloutsos/Pavlo CMU SCS 15-415/615 28

CMU SCS

Query Plan Example

```

SELECT cname, amt
FROM customer, account
WHERE customer.acctno =
      account.acctno
      AND
      account.amt > 1000
    
```

Faloutsos/Pavlo CMU SCS 15-415/615 29

CMU SCS

Today's Class

- Catalog (12.1)
- Intro to Operator Evaluation (12.2,3)
- Typical Query Optimizer (12.6)
- Projection/Aggregation: Sort vs. Hash (14.3.2)

Faloutsos/Pavlo CMU SCS 15-415/615 30

CMU SCS

Duplicate Elimination

```

SELECT DISTINCT bname
FROM account
WHERE amt > 1000
    
```

- What does it do, in English?
- How to execute it?

$\pi_{\text{DISTINCT bname}}(\sigma_{\text{amt} > 1000}(\text{account}))$

Not technically correct because RA doesn't have "DISTINCT"

Faloutsos/Pavlo CMU SCS 15-415/615 31

CMU SCS

Duplicate Elimination

```

SELECT DISTINCT bname
FROM account
WHERE amt > 1000
    
```

Two Choices:

- Sorting
- Hashing

Faloutsos/Pavlo CMU SCS 15-415/615 32

CMU SCS

Sorting Projection

π DISTINCT bname
 σ amt > 1000
 ACCOUNT

acctno	bname	amt
A-123	Redwood	1800
A-789	Downtown	2000
A-123	Perry	1500
A-456	Downtown	1300

Filter

acctno	bname	amt
A-123	Redwood	1800
A-789	Downtown	2000
A-123	Perry	1500
A-456	Downtown	1300

Remove Columns

bname
Redwood
Downtown
Perry
Downtown

Sort

bname
Downtown
Downtown
Perry
Redwood

Eliminate Dupes

Faloutsos/Pavlo CMU SCS 15-415/615 33

CMU SCS

Alternative to Sorting: Hashing!

- What if we don't need the *order* of the sorted data?
 - Forming groups in **GROUP BY**
 - Removing duplicates in **DISTINCT**
- Hashing does this!
 - And may be cheaper than sorting! (why?)
 - But what if table doesn't fit in memory?

Faloutsos/Pavlo CMU SCS 15-415/615 34

CMU SCS

Hashing Projection

- Populate an ephemeral hash table as we iterate over a table.
- For each record, check whether there is already an entry in the hash table:
 - **DISTINCT**: Discard duplicate.
 - **GROUP BY**: Perform aggregate computation.
- Two phase approach.

Faloutsos/Pavlo CMU SCS 15-415/615 35

CMU SCS

Phase 1: Partition

- Use a hash function h_1 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.
- Assume that we have B buffers.

Faloutsos/Pavlo CMU SCS 15-415/615 36

CMU SCS

Phase 1: Partition

π DISTINCT bname
 σ amt > 1000
ACCOUNT

acctno	bname	amt
A-123	Redwood	1800
A-789	Downtown	2000
A-123	Perry	1500
A-456	Downtown	1300

Filter →

acctno	bname	amt
A-123	Redwood	1800
A-789	Downtown	2000
A-123	Perry	1500
A-456	Downtown	1300

Remove Columns →

bname
Redwood
Downtown
Perry
Downtown

Hash (h_1) → **B-1 partitions**

- Redwood
- Downtown
- Downtown
- ⋮
- Perry

Faloutsos/Pavlo CMU SCS 15-415/615 37

CMU SCS

Phase 2: ReHash

- For each partition on disk:
 - Read it into memory and build an in-memory hash table based on a hash function h_2
 - Then go through each bucket of this hash table to bring together matching tuples
- This assumes that each partition fits in memory.

Faloutsos/Pavlo CMU SCS 15-415/615 38

CMU SCS

Phase 2: ReHash

π DISTINCT bname
 σ amt > 1000
ACCOUNT

acctno	bname	amt
A-123	Redwood	1800
A-789	Downtown	2000
A-123	Perry	1500
A-456	Downtown	1300

Partitions From Phase 1

- Redwood → h_2
- Downtown → h_2
- Downtown → h_2
- ⋮
- Perry → h_2

Hash Table

key	value
XXX	Downtown
YYY	Redwood
ZZZ	Perry

Eliminate Dupes

bname
Downtown
Perry
Redwood

Faloutsos/Pavlo CMU SCS 15-415/615 39

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

Faloutsos/Pavlo CMU SCS 15-415/615 40

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: $B \cdot (B-1)$.
 - A table of N blocks needs about $\text{sqrt}(N)$ buffers
 - What assumption do we make?

Faloutsos/Pavlo CMU SCS 15-415/615 41

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: $B \cdot (B-1)$.
 - A table of N blocks needs about $\text{sqrt}(N)$ buffers
 - Assumes hash distributes records evenly!
 - Use a “fudge factor” $f > 1$ for that: we need
 - $B \sim \text{sqrt}(f \cdot N)$

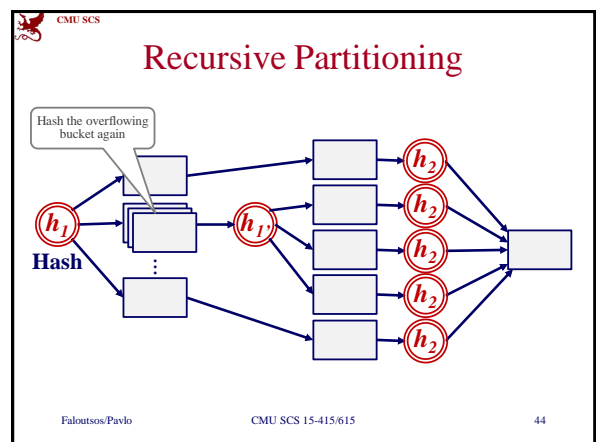
Faloutsos/Pavlo CMU SCS 15-415/615 42

CMU SCS

Analysis

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

Faloutsos/Pavlo CMU SCS 15-415/615 43



CMU SCS

Hashing vs. Sorting

- Which one needs more buffers?

Faloutsos/Pavlo CMU SCS 15-415/615 45

CMU SCS

Hashing vs. Sorting

- **Recall: We can hash a table of size N blocks in \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 \sqrt{N}
 - **Same as hashing!**

Faloutsos/Pavlo CMU SCS 15-415/615 46

CMU SCS

Hashing vs. Sorting

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

Faloutsos/Pavlo CMU SCS 15-415/615 47

CMU SCS

Hashing vs. Sorting

- Choice of **sorting** vs. **hashing** is subtle and depends on optimizations done in each case
- Another optimization when using **sorting** for aggregation:
 - “Early aggregation” of records in sorted runs
- Let’s look at some optimizations for hashing next...

Faloutsos/Pavlo CMU SCS 15-415/615 48

CMU SCS

Hashing: We Can Do Better!

- Combine the summarization into the hashing process - How?

Faloutsos/Pavlo CMU SCS 15-415/615 49

CMU SCS

Hashing: We Can Do Better!

- During the ReHash phase, store pairs of the form **<GroupKey, RunningVal>**
- When we want to insert a new tuple into the hash table:
 - If we find a matching **GroupKey**, just update the **RunningVal** appropriately
 - Else insert a new **<GroupKey, RunningVal>**

Faloutsos/Pavlo CMU SCS 15-415/615 50

CMU SCS

Hashing Aggregation

```
SELECT bname, SUM(amt)
FROM account
GROUP BY bname
```

key	value
XXX	<Redwood, 1200->
YYY	<Downtown, 1000->
ZZZ	<Perry, 1500->

bname	SUM(amt)
Redwood	4355
Downtown	6895
Perry	7901

Faloutsos/Pavlo CMU SCS 15-415/615 51

CMU SCS

Hashing Aggregation

- What's the benefit?
- How many entries will we have to handle?
 - Number of distinct values of GroupKeys columns
 - Not the number of tuples!!
 - Also probably "narrower" than the tuples

Faloutsos/Pavlo CMU SCS 15-415/615 52

CMU SCS

So, hashing is better...right?

- Any caveats?

Faloutsos/Pavlo CMU SCS 15-415/615 53

CMU SCS

So, hashing is 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

Faloutsos/Pavlo CMU SCS 15-415/615 54

CMU SCS

Summary

- Query processing architecture:
 - Query optimizer translates SQL to a query plan
 - Query executor “interprets” the plan
- **Hashing** is a useful alternative to **sorting** for duplicate elimination / group-by
 - Both are valuable techniques for a DBMS

Faloutsos/Pavlo CMU SCS 15-415/615 55

CMU SCS

Next Class

- How to actually use indexes.
- ➔ • Join algorithms.
- More query optimization.

Faloutsos/Pavlo CMU SCS 15-415/615 56