

EECS 339: Introduction to Database Systems

Midterm Exam
February 17, 2015

Name:

I Schema Design

1. **5 points:** You are designing a real estate database with the following requirements:

- Properties describe houses or plots of land. They have an id, address, an area in square feet, a number of bedrooms, a count of bathrooms, an asking price, and the date it was listed
- A house has an ownership history, where each homeowner has an id, name, date of birth, a social security number, year they bought it, end ownership date, and price paid. An owner may own any number of properties, either in the past or concurrently.
- Properties each have a category, such as commercial, residential, or land. Each distinct category has an id.

Draw an entity relationship diagram for this database. Please draw entities as squares, attributes as ovals, and denote relationships as diamonds between pairs of entities, with a label of the form “1:1”, “N:1” or “1:N”, where “1:N” indicates that a single entity on the left side has a relationship with many entities on the right, but each entity on the right has a relationship with only one entity on the left. Give each entity, relationship, and attribute a name.

2. **5 points:** Using your entity-relationship diagram, create a schema for this database. Use the following format:

`tablename (attribute1, attribute2, attribute3)`

Underline the one or more attributes for each primary key.

3. 5 points: Is this schema in Boyce-Codd Normal Form?

YES NO

Why or why not? If it has any redundancies or anomalies explain them below.

II Query Writing

For the schema below, express the following questions as SQL queries.

4. 15 points: A university's database has the schema:

```
students (sid, s_name, dob)
courses (cid, c_name, year, quarter, meeting_time, building, room)
grades (s_sid, c_cid, grade)
```

(a) In what building did the class with id 'eecs339' meet during the winter quarter of 2015?

(b) Create a list of each student's id, name, and their GPA.

(c) List all course ids, course names, and the number of enrolled students. Order the output alphabetically by name.

(d) List the names of all students who got a grade of 4.0 in the course with id 'eecs339'.

- (e) What is the `cid` of the class(es) where the students had the highest average grade?

III Query Rewriting

5. **10 points:** For the following queries, can they be rewritten for faster query execution? If so, rewrite the query. If not, explain why not.

(a)

```
SELECT *
FROM EMPLOYEES
WHERE salary * (1 + 0.05) > 100000;}
```

Rewrite possible?

YES **NO**

Explanation:

(b)

```
WITH ninjas as (SELECT *
FROM employees as e, departments as d
WHERE e.dept_no = d.dept_no and d.dept_name = 'ninjas')
SELECT n.name
FROM ninjas as n, favorite_colors as fc
WHERE n.emp_id = fc.emp_id and fc.color = 'black';
```

Rewrite possible?

YES **NO**

Explanation:

(c)

```
with feedingCount as (
SELECT aid, count(*) as cnt
FROM feedings
GROUP BY aid)
SELECT a.aid, a.name, f.cnt
FROM animals as a, feedingCount as f
WHERE f.aid = a.aid AND f.cnt > 2;
```

Rewrite possible?

YES **NO**

Explanation:

```
(d)  with bigTransactions as
      (SELECT  distinct acct_id
       FROM transactions
       WHERE amt > 10000)
SELECT  a.name, a.addr
FROM accounts
WHERE EXISTS (
SELECT * FROM bigTransactions
WHERE bigTransactions.acct_id = accounts.acct_id);
```

Rewrite possible?

YES **NO**

Explanation:

IV Short Answer

IV.1 Memory Management

- 6. 5 points:** Based on the DBMIN algorithm from the paper, "An Evaluation of Buffer Management Strategies for Relational Database Systems", what percent of buffer pool accesses will be hits for each of the following situations create for LRU eviction policies? Assume the system has up to 50 MB of buffer pool space. Each database page is 5 MB.

(a) A 100 MB table with 10 tuples per page that is sequentially scanned.

(b) A nested loop join with the outer relation having 100 pages and 1,000 tuples total, and the inner relation having 5 pages with 1 tuple per page.

- 7. 5 points:** You have a queue of queries with the following working set sizes (in pages):

13, 16, 11, 18, 20, 5, 2, 15, 7, 3, 10, 1, 8, 14, 4

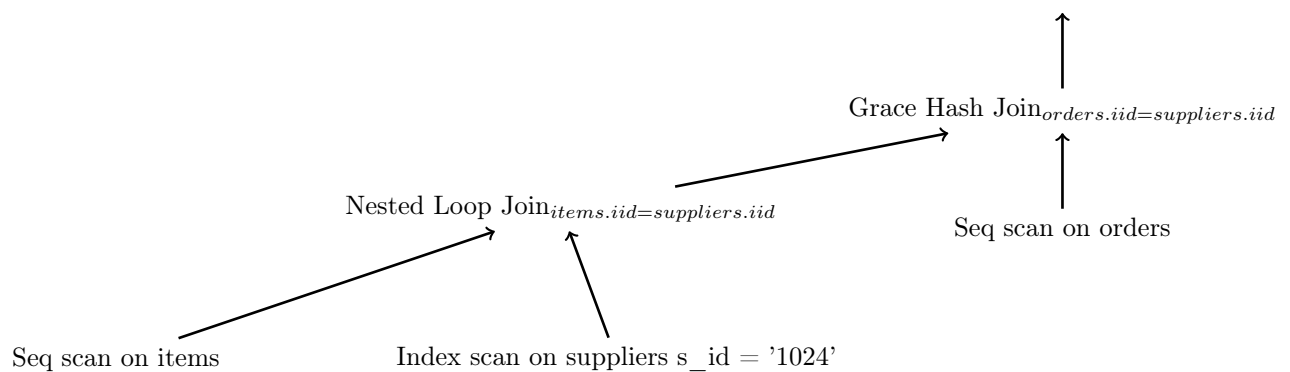
Your buffer pool contains 30 pages. If the queries all have the same duration, in what order will they be executed?

IV.2 Join Algorithms

8. 5 points: For the query:

```
SELECT *
FROM items, suppliers, orders
WHERE items.iid = suppliers.iid AND suppliers.iid = orders.iid
AND suppliers.s_id = 1024;
```

The query optimizer produces the following query plan:



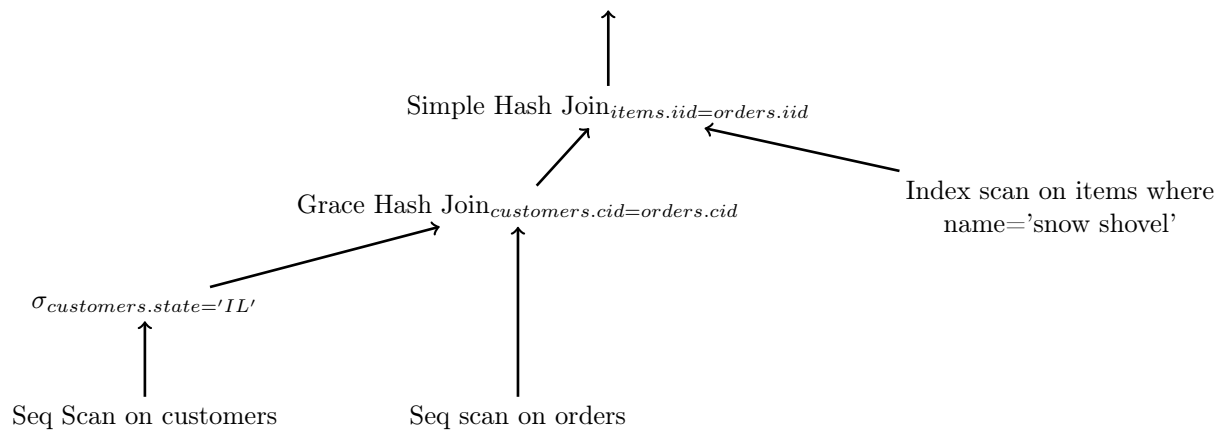
(a) Why are we using a nested loop join in $\text{items} \bowtie \text{suppliers}$?

(b) Explain the use of a grace hash join in the second join.

9. 5 points: For the query:

```
SELECT *
FROM customers, orders, items
WHERE customers.cid = orders.cid AND items.iid = orders.iid
AND customers.state = 'IL' and items.name = 'snow shovel';
```

The query optimizer produces the following query plan:



- (a) Why are we using a grace hash join in $\text{customers} \bowtie \text{orders}$?
 (b) Justify the use of a simple hash join in the second join.

V Indexing and Access Paths

10. 5 points: We are querying the employee table of a business with the following schema:

`employee (empl_id, name, salary, dept_no, building_no, managed_by)`

For the following scenarios, which index will be the fastest?

1. B+ tree
2. Extendible hash
3. Bitmap index
4. Sequential Scan

The company has three buildings, and each has an equal number of employees. The table is much bigger than memory and clustered on its primary key, employee id. It has 1000 managers, each having a different number of direct reports in `managed_by`.

Match each of the letters above with one scenario below:

- (a) Count the number of employees in the Glass building.
 (b) Sum up the total salaries paid over all employees.
 (c) Look up specific employees one at a time by id to give them raises.
 (d) List the employees managed by John.

VI Query Planning

Consider the following query:

```
SELECT *  
FROM A, B, C  
WHERE A.v = B.v AND A.w = C.w  
AND A.x > 7 and B.v = 5  
AND C.w <= 3;
```

Assume that all of the attributes have a range of 1...10, inclusive. A has 100 tuples, B has 200, and C has 500.

11. **5 points:** Estimate the number of tuples that would be initially selected from each of the three relations if all of the non-join predicates are applied to them before any join processing begins.
12. **5 points:** If we join in the order $A \bowtie B \bowtie C$, what is the expected output cardinality? Show your arithmetic.
13. **5 points:** Draw a query plan tree for $A \bowtie B \bowtie C$. Label the estimated size of each intermediate result and the output cardinality.

VII Query Optimization

Consider a query of the form $A \bowtie B \bowtie C \bowtie D \bowtie E$.

14. **5 points:** With this five-way join, how many join orderings are possible if we enumerate all possible orderings? Do not consider different parenthesizations, i.e., only consider left-deep plan structures.
15. **8 points:** What fraction of these plan evaluations can we eliminate using dynamic programming? In other words, if we generate the plan incrementally rather than considering whole, 5-relation ones, what is the ratio of partial orderings to whole ones? Evaluate orderings like $A \bowtie B$ and $B \bowtie A$ as a single plan.

- 16. 7 points:** What if the query is $A \bowtie_{A.v=B.v} B \bowtie_{B.v=C.v} C \bowtie_{C.w=D.w} D \bowtie_{D.w=E.w} E$? What partial plans from the dynamic programming approach can we rule out by not considering cross joins? Again, consider $A \bowtie B$ and $B \bowtie A$ the same for this analysis. List the eliminated sub-plans below.