# COMP7640
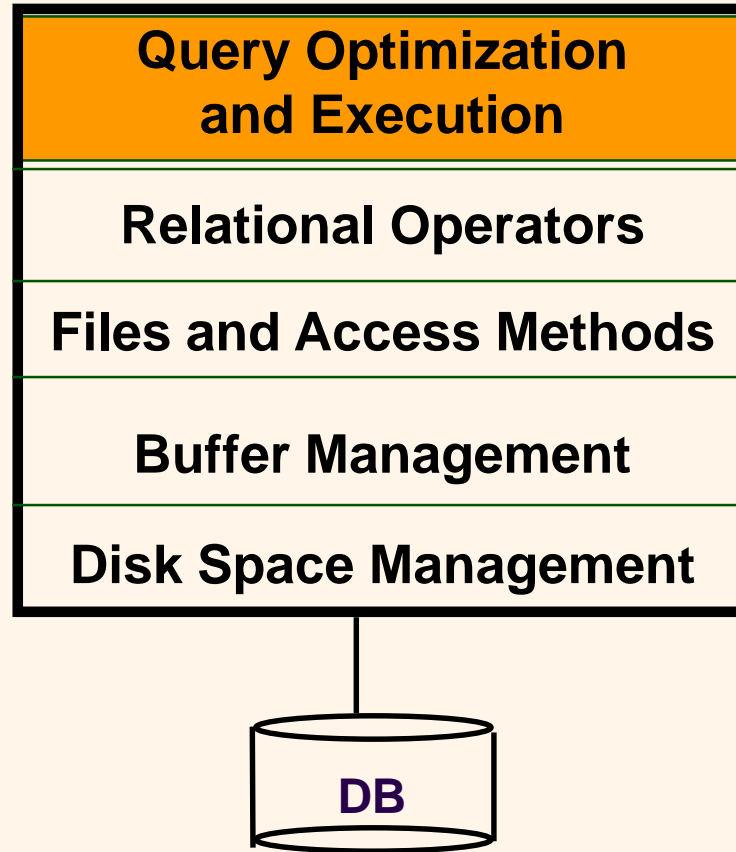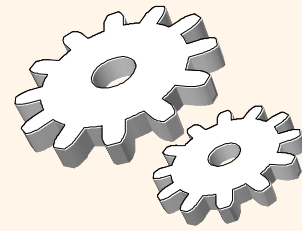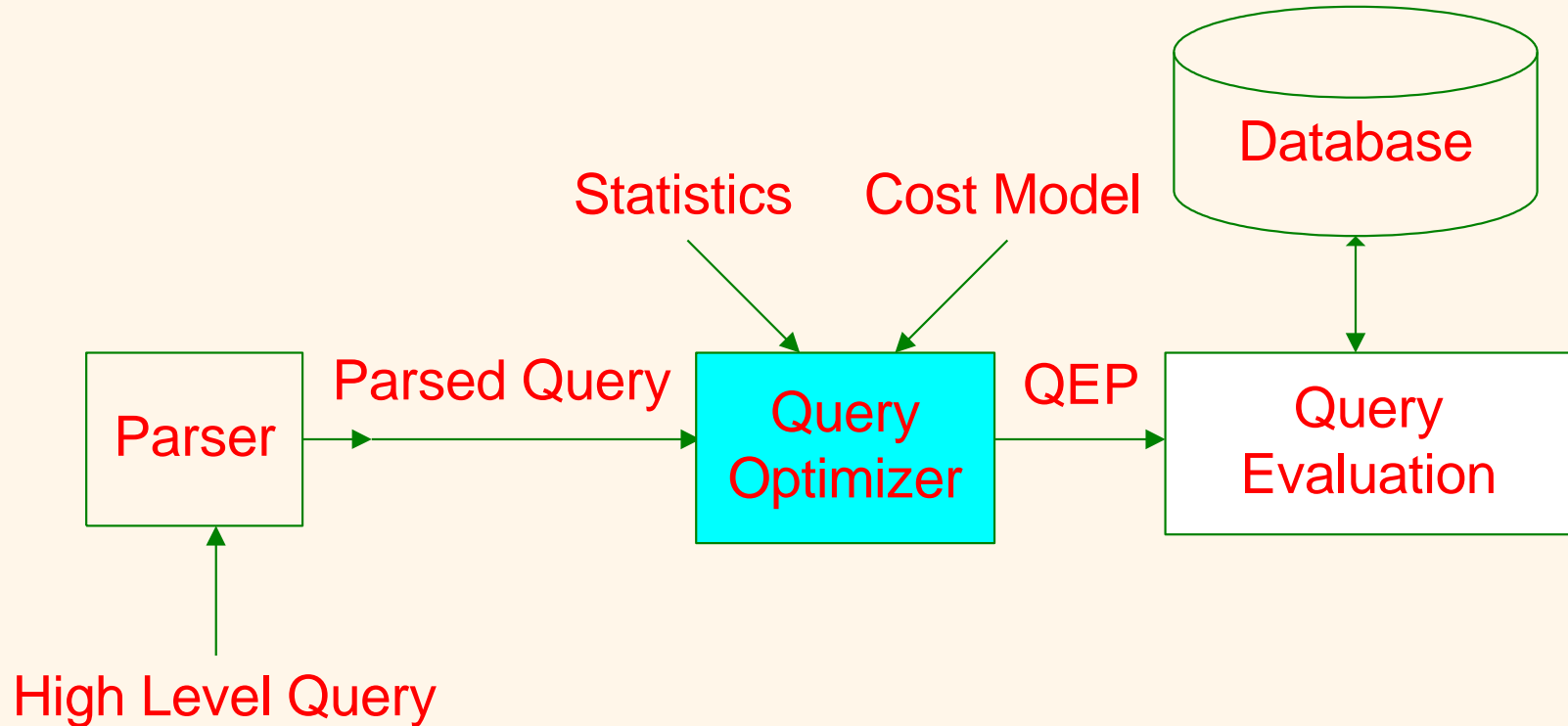# Database Systems & Administration

*Query Optimization*

# *Where Are We Now?*

| Query Optimization and Execution |
|:---:|
| **Relational Operators** |
| **Files and Access Methods** |
| **Buffer Management** |
| **Disk Space Management** |

**DB**

# *Processing a High-Level Query*

Statistics   Cost Model

Database

Parser → **Parsed Query** → Query Optimizer → **QEP** → Query Evaluation

High Level Query

# *In Query Evaluation*
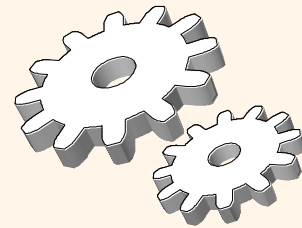
❖ Various access paths for relational operators
  ▪ Selection
    • Sorted file
    • Index (B+ tree/ hash index)
  ▪ Projection
    • Sort-based projection
  ▪ Join
    • Simple nested-loop join
    • Page-oriented nested-loop join
❖ Only evaluate the cost of a *single* relational operator

# SQL Queries In Practice

**SELECT** S.sname
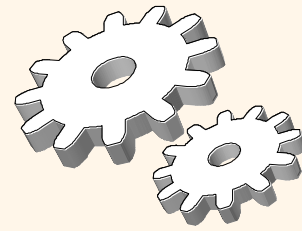
**FROM** Students S, CourseEnrolled E

**WHERE**  S.sid= E.sid **AND** E.cid = 3220 **AND** S.gpa > 3.0

$$\pi_{\text{sname}}\left(\sigma_{\text{gpa}>3.0\,\wedge\,\text{cid}=3220}(S \bowtie_{\text{S.sid=E.sid}} E)\right)$$

- ❖ A query is basically a *relational algebra expression* (a set of ordered *relational operators*)

- ❖ A query can be represented by *multiple* *relational algebra expressions* (*relational operators* in different orders)

- ❖ A query can involve *multiple* *relational operators* in its *relational algebra expression*

- ❖ Each relational operator can be implemented via *multiple access paths*

# *Query Evaluation Plan (QEP)*

❖ A *query evaluation plan* (or *query execution plan,* QEP) tells the DBMS how to execute the SQL query. It specifies

  ▪ A *relational algebra expression*

    • What operations we need to execute

    • What execution orders of these operations are

  ▪ Access paths for each *relational operator* in the *relational algebra expression*

# *Query Optimization*

❖ The goal of *query optimization* is to find the QEP with *lest* I/O cost before execution

❖ For a given query, how to choose a *good* *query evaluation plan* (or *query execution plan*, QEP)
  - Enumerate *alternative* QEPs
  - Estimate *cost* (I/O cost) of *each* enumerated QEP
  - Choose the QEP with *least* cost

# *How to Enumerate Alternative QEPs?*
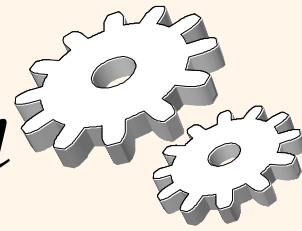
❖ Transform the *relational algebra expression* to its equivalent forms (with different orders of operators)

- Equivalent rules (ensure that the results of the alternative plan are correct)
- Different (equivalent) expressions can significantly affect the I/O cost

❖ Choose *appropriate access paths* for each *relational operator* in the *relational algebra expression*

# *Decomposition Rule for Selection Operations*

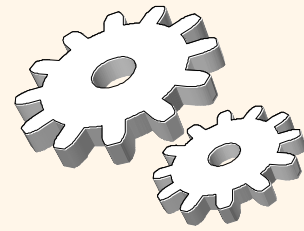❖ *Decompose* the selection operation to *multiple* selection operations

❖ Consider the selection operation with conditions $\theta_1$ and $\theta_2$ for the relation $R$. We have:
$$\sigma_{\theta_1 \wedge \theta_2}(R) = \sigma_{\theta_1}(\sigma_{\theta_2}(R))$$

❖ **Example**: Let the relation be $R(a, b, c)$. We have:
$$\sigma_{a>20 \,\wedge\, b \leq 20}(R) = \sigma_{a>20}(\sigma_{b \leq 20}(R))$$

# *Commutative Rule for Selection Operations*

❖ *Reverse* the order of *consecutive* selection operations
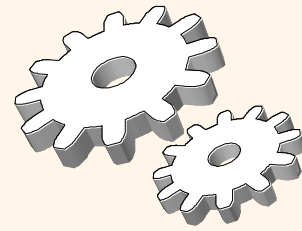
❖ Consider the selection operation with conditions $\theta_1$ and $\theta_2$ for the relation $R$. We have:

$$\sigma_{\theta_1}\left(\sigma_{\theta_2}(R)\right) = \sigma_{\theta_2}\left(\sigma_{\theta_1}(R)\right)$$

❖ **Example**: Let the relation be $R(a, b, c)$. We have:

$$\sigma_{a>20}\left(\sigma_{b\leq20}(R)\right) = \sigma_{b\leq20}\left(\sigma_{a>20}(R)\right)$$

# *Omission Rule* *for Projection Operations*

❖ Given *multiple* projection operations, only the *final* one should be retained.

❖ Consider the projection operations with the sets of attributes $L_1, L_2, \ldots, L_n$, where $L_1 \subseteq L_2 \subseteq \cdots \subseteq L_n$, for the relation $R$.

$$\pi_{L_1}\left(\pi_{L_2}\left(\cdots\left(\pi_{L_n}(R)\right)\cdots\right)\right) = \pi_{L_1}(R)$$

**$R$**

| a | b | c | d |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 1 | 5 | 0 | 0 |
| 2 | 5 | 0 | 0 |
| 3 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 |

$\pi_{a,b,c}(R)$ →

**$R_1$**

| a | b | c |
|---|---|---|
| 0 | 0 | 0 |
| 1 | 5 | 0 |
| 2 | 5 | 0 |
| 3 | 0 | 0 |
| 4 | 0 | 0 |

$\pi_{a,b}(R_1)$ →

**$R_2$**

| a | b |
|---|---|
| 0 | 0 |
| 1 | 5 |
| 2 | 5 |
| 3 | 0 |
| 4 | 0 |

$\pi_a(R_2)$ →

**$R_3$**

| a |
|---|
| 0 |
| 1 |
| 2 |
| 3 |
| 4 |

# *Commutative Rule* for *Selection and Projection Operations*

- ❖ *Reverse* the order of *consecutive* <u>selection</u> and <u>projection</u> operations with specific condition

- ❖ Consider the <u>projection</u> operation with the set of attributes $L$ and the <u>selection</u> operation with condition $\theta$. If the *condition $\theta$ only involves those attributes in L*. We have:

$$\pi_L(\sigma_\theta(R)) = \sigma_\theta(\pi_L(R))$$

# *Commutative Rule for Selection and Projection Operations*

$$\theta = b > 0$$
$$L = \{a, b\}$$

$R$

| a | b | c | d |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 1 | 5 | 0 | 0 |
| 2 | 5 | 0 | 0 |
| 3 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 |

$\sigma_\theta(R)$

| a | b | c | d |
|---|---|---|---|
| 1 | 5 | 0 | 0 |
| 2 | 5 | 0 | 0 |

$\pi_L(R)$

| a | b |
|---|---|
| 1 | 5 |
| 2 | 5 |

$\pi_L(R)$

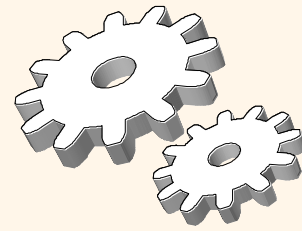| a | b |
|---|---|
| 0 | 0 |
| 1 | 5 |
| 2 | 5 |
| 3 | 0 |
| 4 | 0 |

$\sigma_\theta(R)$

❖ **Example**: Let the relation be $R(a, b, c)$. We have:
$$\pi_{a,b}(\sigma_{a>20}(R)) = \sigma_{a>20}(\pi_{a,b}(R))$$
$$\pi_{a,b}(\sigma_{a>10 \,\wedge\, b\leq30}(R)) = \sigma_{a>10 \,\wedge\, b\leq30}(\pi_{a,b}(R))$$
$$\pi_{a,b}(\sigma_{c>10 \,\wedge\, b\leq30}(R)) \neq \sigma_{c>10 \,\wedge\, b\leq30}(\pi_{a,b}(R))$$

13

# *Distributive Rule 1 for Selection and Join Operations*

❖ *Distribute* the <u>selection</u> operation in the <u>join</u> operation with specific condition

❖ Consider the <u>selection</u> operation with condition $\theta_1$ and the <u>join</u> operation with condition $\theta$ for two relations $R$ and $S$. If *all the attributes in $\theta_1$ <u>only</u> involve the attributes of the relation $R$.*

$$\sigma_{\theta_1}(R \bowtie_\theta S) = \sigma_{\theta_1}(R) \bowtie_\theta S$$

# *Distributive Rule 1 for Selection and Join Operations*

❖ **Example**:

*Table S*

| sid | sname | gpa | age |
|-----|-------|-----|-----|
| 22 | simon | 3.6 | 20 |
| 31 | kelvin | 3.5 | 21 |
| 58 | karen | 3.5 | 18 |

*Table E*

| sid | cid | day |
|-----|-----|-----|
| 22 | 2440 | 10/01/04 |
| 22 | 3220 | 10/12/03 |
| 58 | 3820 | 11/01/04 |

$$\sigma_{S.gpa > 3.5}(S \bowtie_{S.sid=E.sid} E) \quad \Leftrightarrow \quad \sigma_{S.gpa > 3.5}(S) \bowtie_{S.sid=E.sid} E$$

$$\sigma_{E.cid > 2440}(S \bowtie_{S.sid=E.sid} E) \quad \not\Leftrightarrow \quad \boxed{\sigma_{E.cid > 2440}(S)} \bowtie_{S.sid=E.sid} E$$

❌
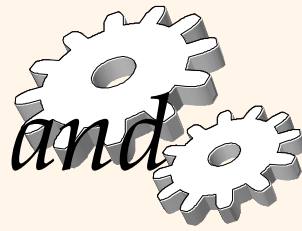
# *Distributive Rule* 2 *for Selection and Join Operations*

❖ *Distribute* the <u>selection</u> operation in the <u>join</u> operation with specific condition

❖ Consider the <u>selection</u> operation with conditions $\theta_1$ and $\theta_2$ and the <u>join</u> operation with condition $\theta$ for two relations $R$ and $S$. If *all the attributes in $\theta_1$ and $\theta_2$ <u>only</u> involve the attributes of the relations $R$ and $S$,* respectively.

$$\sigma_{\theta_1 \wedge \theta_2}(R \bowtie_\theta S) = \sigma_{\theta_1}(R) \bowtie_\theta \sigma_{\theta_2}(S)$$

# *Distributive Rule 2 for Selection and Join Operations*

❖ **Example**:

*Table S*

| sid | sname | gpa | age |
|-----|-------|-----|-----|
| 22  | simon | 3.6 | 20  |
| 31  | kelvin | 3.5 | 21  |
| 58  | karen | 3.5 | 18  |

*Table E*

| sid | cid | day |
|-----|-----|-----|
| 22  | 2440 | 10/01/04 |
| 22  | 3220 | 10/12/03 |
| 58  | 3820 | 11/01/04 |

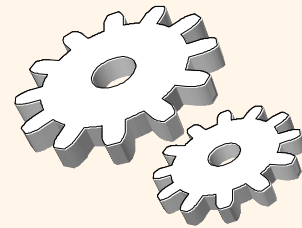$$\sigma_{S.gpa>3.5 \wedge E.cid=2440}(S \bowtie_{S.sid=E.sid} E)$$

$$\Updownarrow$$

$$\sigma_{S.gpa>3.5}(S) \bowtie_{S.sid=E.sid} \sigma_{E.cid=2440}(E)$$

# *Summary*

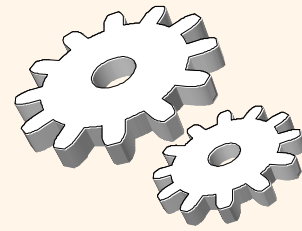| Relational Operator | Rules |
|---|---|
| Selection $\sigma$ | $\sigma_{\theta_1 \wedge \theta_2}(R) \Leftrightarrow \sigma_{\theta_1}(\sigma_{\theta_2}(R))$ **Decomposition Rule** |
| | $\sigma_{\theta_1}(\sigma_{\theta_2}(R)) \Leftrightarrow \sigma_{\theta_2}(\sigma_{\theta_1}(R))$ **Commutative Rule** |
| Projection $\pi$ | $\pi_{L_1}\left(\pi_{L_2}\left(\cdots\left(\pi_{L_n}(R)\right)\cdots\right)\right) \Leftrightarrow \pi_{L_1}(R), L_1 \subseteq L_2 \subseteq \cdots \subseteq L_n$ **Omission Rule** |
| Selection & Projection $\sigma$ & $\pi$ | $\pi_L(\sigma_\theta(R)) \Leftrightarrow \sigma_\theta(\pi_L(R)), \theta \subseteq L$ **Commutative Rule** |
| Selection & Join $\sigma$ & $\bowtie$ | $\sigma_{\theta_1}(R \bowtie_\theta S) \Leftrightarrow \sigma_{\theta_1}(R) \bowtie_\theta S, \theta_1 \subseteq R$ **Distributive Rule 1** |
| | $\sigma_{\theta_1 \wedge \theta_2}(R \bowtie_\theta S) \Leftrightarrow \sigma_{\theta_1}(R) \bowtie_\theta \sigma_{\theta_2}(S), \theta_1 \subseteq R, \theta_2 \subseteq S$ **Distributive Rule 2** |

# *Question 1*

❖ Given a relational algebra expression

$$\pi_{\text{sname,gpa,cid}} \left( \sigma_{\text{gpa}>3.0 \, \wedge \, \text{cid}=3220}(S \bowtie_{\text{S.sid}=\text{E.sid}} E) \right)$$

List its *two* equivalent relational algebra expressions.

# *How to Enumerate Alternative QEPs?*

❖ Enumerate *alternative relational algebra expressions* based on the rules (Slides 9-18)

❖ Specify *access paths* for each *relational operator* (Lecture 10: Query Evaluation)

❖ For each enumerated QEP, we represent it by
  ▪ The *relational algebra tree* of its *relational algebra expression*
  ▪ Annotating at each node to indicate the *access path* for each *relational operator*

# *Example Instances*

**Students S (sid: *integer*, sname: *string*, gpa: *real*, age: *integer*)**
**CourseEnrolled E (sid: *integer*, cid: *string*, day: *date*)**

### *Table S*

| sid | sname | gpa | age |
|-----|-------|-----|-----|
| 22 | simon | 3.6 | 20 |
| 31 | kelvin | 3.5 | 21 |
| 58 | karen | 3.5 | 18 |

*# of records: 100,000*
*# of pages: 1,000*

### *Table E*

| sid | cid | day |
|-----|-----|-----|
| 22 | 2440 | 10/01/04 |
| 22 | 3220 | 10/12/03 |
| 58 | 3820 | 11/01/04 |

*# of records: 200,000*
*# of pages: 500*

**Assumption 1**: The gpa attribute follows the *uniform distribution* in the range 0 to 4.
**Assumption 2**: There are 50 courses (i.e., 50 cids) in the Table E and these cids are *uniformly distributed* in this relation.

# *Example*

❖ Example SQL query

SELECT S.sname

FROM Students S, CourseEnrolled E

WHERE  S.sid= E.sid **AND** E.cid = 3220 **AND** S.gpa > 3.0

❖ Relational algebra expression

$$\pi_{\text{sname}}\left(\sigma_{\text{gpa}>3.0 \,\wedge\, \text{cid}=3220}(S \bowtie_{\text{S.sid=E.sid}} E)\right)$$
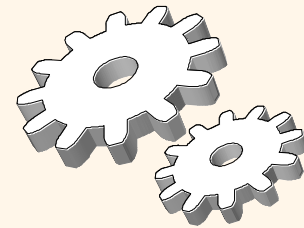
# *Query Evaluation Plan*

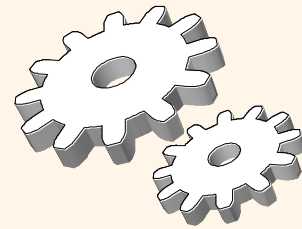❖ Represent this *relational algebra expression* as a *relational algebra tree* with *access paths*

$\pi_{\text{sname}}$      **(On-the-fly)**

**done in the memory**

$\sigma_{\text{cid}=3220}$ $\wedge$ $\sigma_{\text{gpa}>3.0}$      **(On-the-fly)**

**What are they? (Pipelined Evaluation)**

$\bowtie_{\text{sid=sid}}$      **(Page-oriented nested-loop join)**

**Access path for join operation**

| sid | sname | gpa | age |
|-----|-------|-----|-----|
| 22  | simon | 3.6 | 20  |
| 31  | kelvin | 3.5 | 21 |
| 58  | karen | 3.5 | 18  |

**Students**

**CourseEnrolled**

| sid | cid | day |
|-----|-----|-----|
| 22  | 2440 | 10/01/04 |
| 22  | 3220 | 10/12/03 |
| 58  | 3820 | 11/01/04 |

# *Query Evaluation Plan*

## ❖ **Pipelined Evaluation**

- ▪ Result of one operator pipelined to another operator *without creating a temporary relation* to hold intermediate result
  - Each record produced by an operator in the memory will be directly sent to the next operators in the memory without writing it to or reading it from the disk. The subsequent operations are done in the memory and entail no costs.

- ▪ Lower overhead
  - Avoid the cost of writing out intermediate results
  - Avoid reading those results to the main memory

# *Query Evaluation Plan*

## ❖ **Pipelined Evaluation**

# *Query Evaluation Plan*

## ❖ **Pipelined Evaluation**

$\pi_{\text{sname}}$

| sname |
|-------|
| simon |

$\sigma_{\text{cid}=3220} \wedge \sigma_{\text{gpa}>3.0}$

**Memory**

| sid | sname | gpa | age | sid | cid | day |
|-----|-------|-----|-----|-----|------|----------|
| 22 | simon | 3.6 | 20 | 22 | 3220 | 10/12/03 |

✔

| sid | sname | gpa | age | sid | cid | day |
|-----|-------|-----|-----|-----|------|----------|
| 22 | simon | 3.6 | 20 | 22 | 3220 | 10/12/03 |

✔

$\bowtie_{\text{sid=sid}}$

**Read pages**   **Read pages**

**Disk**

| sid | sname | gpa | age |
|-----|-------|-----|-----|
| 22 | simon | 3.6 | 20 |
| 31 | kelvin | 3.5 | 21 |
| 58 | karen | 3.5 | 18 |

**S**

**E**

| sid | cid | day |
|-----|------|----------|
| 22 | 2440 | 10/01/04 |
| 22 | 3220 | 10/12/03 |
| 58 | 3820 | 11/01/04 |

**Disk**

# *Cost of Evaluating the Plan*

❖ <u>No</u> cost for <u>selection</u> and <u>projection</u> operations (pipelined evaluation, done in the memory)

❖ The cost is: 500 + 1000 * 500 = 500,500 I/Os

$\pi_{\text{sname}}$ **(On-the-fly)**
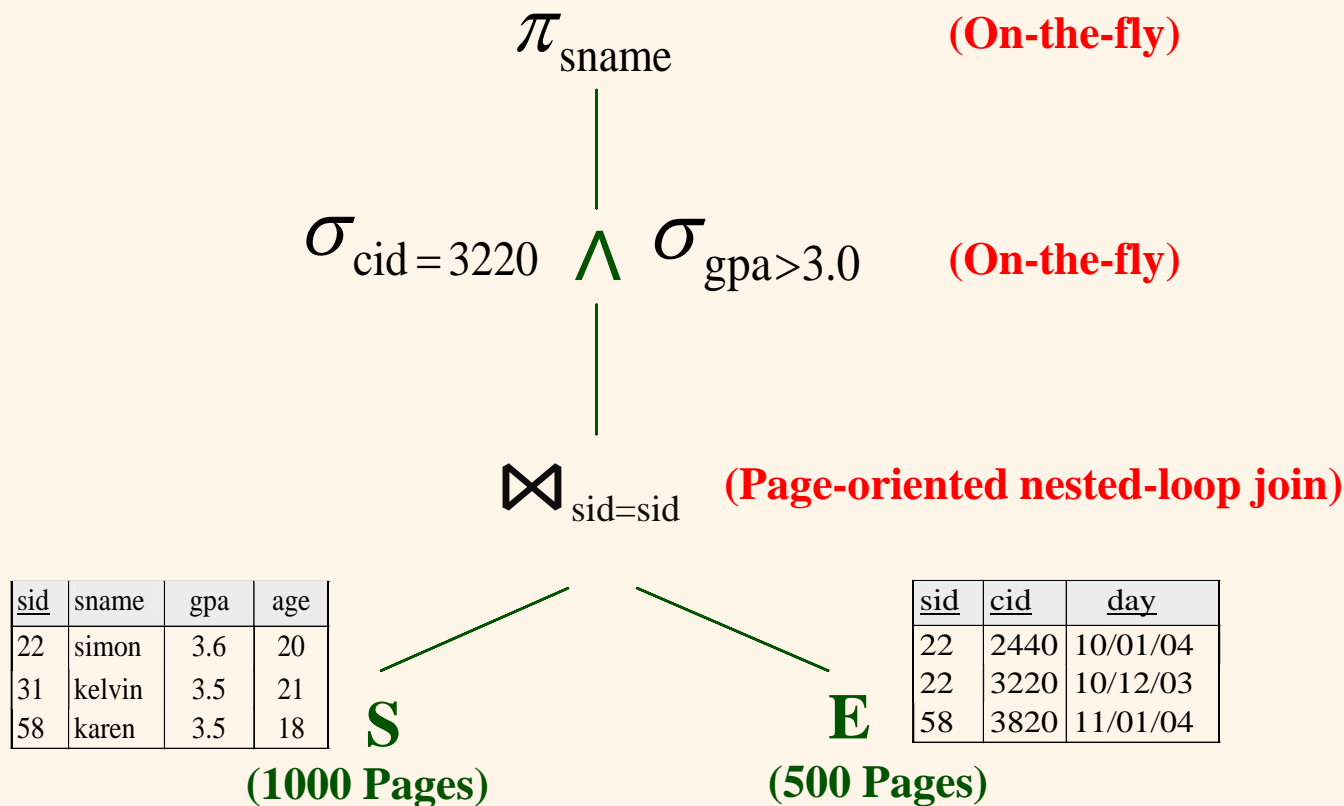
$\sigma_{\text{cid}=3220} \wedge \sigma_{\text{gpa}>3.0}$ **(On-the-fly)**

$\bowtie_{\text{sid=sid}}$ **(Page-oriented nested-loop join)**

| sid | sname | gpa | age |
|-----|-------|-----|-----|
| 22  | simon | 3.6 | 20  |
| 31  | kelvin| 3.5 | 21  |
| 58  | karen | 3.5 | 18  |

**S**
**(1000 Pages)**

| sid | cid  | day      |
|-----|------|----------|
| 22  | 2440 | 10/01/04 |
| 22  | 3220 | 10/12/03 |
| 58  | 3820 | 11/01/04 |

**E**
**(500 Pages)**

# *Alternative Plan 1*

| sid | sname | gpa | age |
|-----|-------|-----|-----|
| 22 | simon | 3.6 | 20 |
| 31 | kelvin | 3.5 | 21 |
| 58 | karen | 3.5 | 18 |

*Table E*

| sid | cid | day |
|-----|-----|-----|
| 22 | 2440 | 10/01/04 |
| 22 | 3220 | 10/12/03 |
| 58 | 3820 | 11/01/04 |

❖ Using distributive rule 2, we can move $\sigma_{gpa>3.0}$ to S and $\sigma_{cid=3220}$ to E
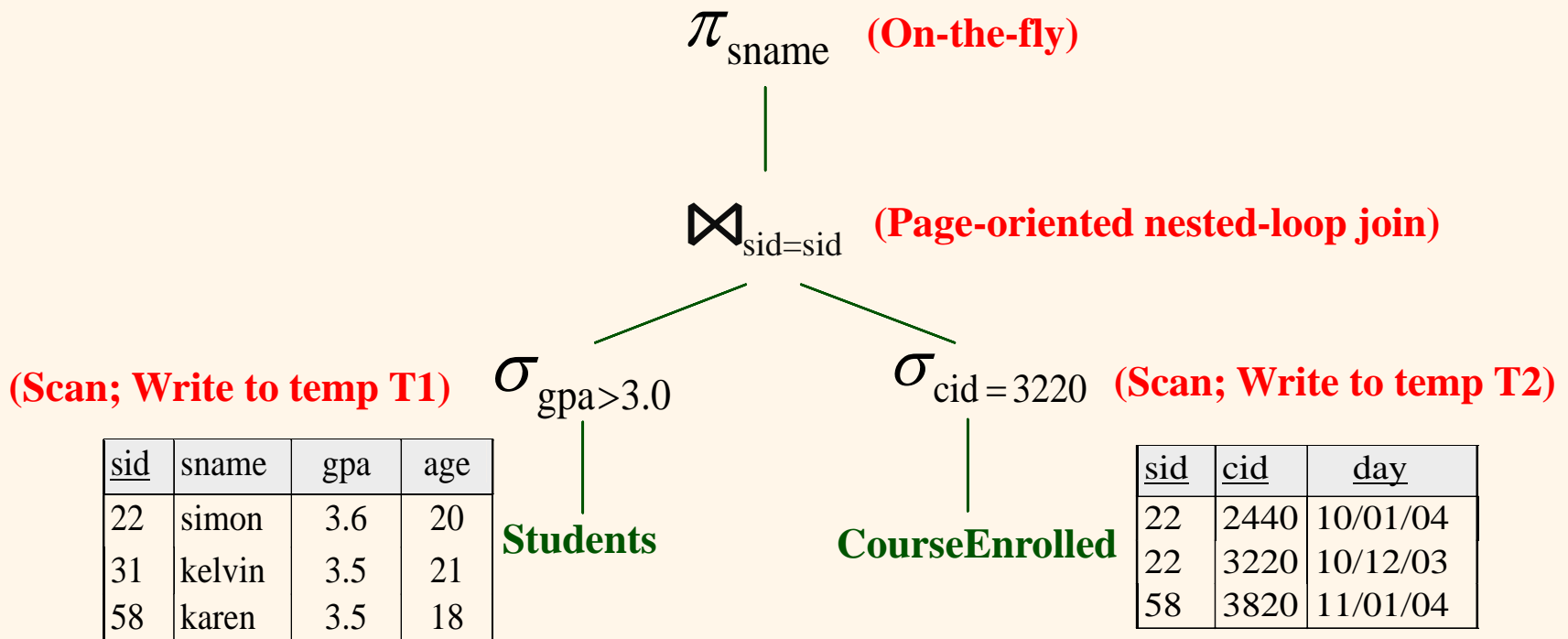
$$\pi_{\text{sname}}\left(\sigma_{\text{gpa>3.0} \wedge \text{cid=3220}}(S \bowtie_{\text{S.sid=E.sid}} E)\right)$$

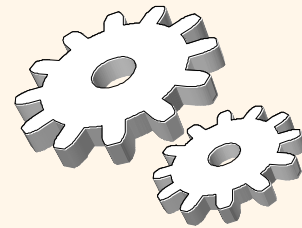$$= \pi_{\text{sname}}\left(\sigma_{\text{gpa>3.0}}(S) \bowtie_{\text{S.sid=E.sid}} \sigma_{\text{cid=3220}}(E)\right)$$

# *Alternative Plan 1*
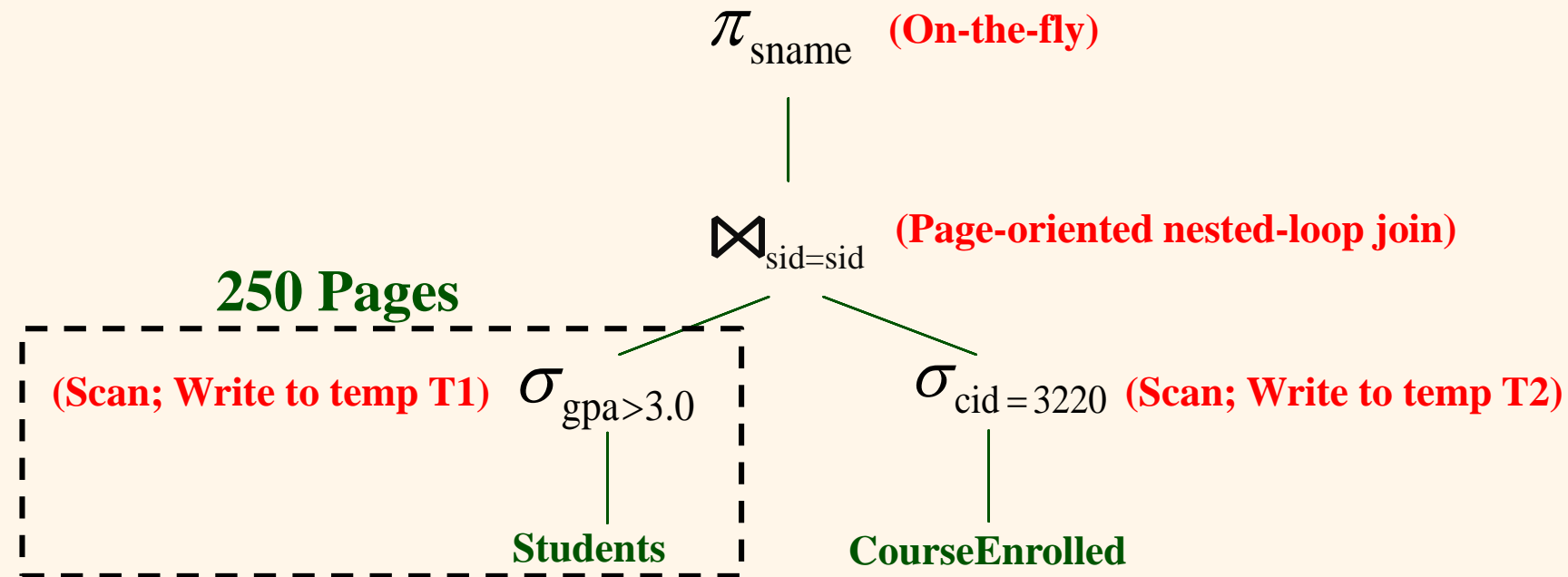
$$\pi_{\text{sname}} \left( \sigma_{\text{gpa}>3.0}(S) \bowtie_{\text{S.sid=E.sid}} \sigma_{\text{cid}=3220}(E) \right)$$

❖ The relational algebra tree is:

$\pi_{\text{sname}}$ **(On-the-fly)**

$\bowtie_{\text{sid=sid}}$ **(Page-oriented nested-loop join)**

**(Scan; Write to temp T1)** $\sigma_{\text{gpa}>3.0}$     $\sigma_{\text{cid}=3220}$ **(Scan; Write to temp T2)**

| sid | sname | gpa | age |
|-----|-------|-----|-----|
| 22  | simon | 3.6 | 20  |
| 31  | kelvin | 3.5 | 21 |
| 58  | karen | 3.5 | 18  |

**Students**

**CourseEnrolled**

| sid | cid  | day      |
|-----|------|----------|
| 22  | 2440 | 10/01/04 |
| 22  | 3220 | 10/12/03 |
| 58  | 3820 | 11/01/04 |

# *Cost of Evaluating the Alternative Plan 1*

$$\pi_{\text{sname}} \quad \text{(On-the-fly)}$$

$$\bowtie_{\text{sid=sid}} \quad \text{(Page-oriented nested-loop join)}$$

**250 Pages**

**(Scan; Write to temp T1)** $\sigma_{\text{gpa}>3.0}$

$\sigma_{\text{cid}=3220}$ **(Scan; Write to temp T2)**

**Students**

**CourseEnrolled**

**Cost of $\sigma_{\text{gpa} > 3.0}$ (S) = Cost to scan S + Cost to write T1**
**= 1000 pages + size(T1)**
**Since the gpa attribute follows the uniform distribution in the range 0 to 4 (Assumption 1), we have:**
**size(T1) = 250 pages**

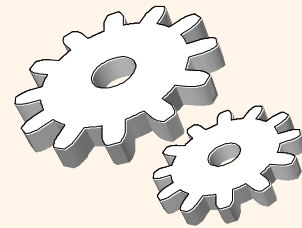# Cost of Evaluating the Alternative Plan 1

$\pi_{\text{sname}}$ **(On-the-fly)**

$\bowtie_{\text{sid=sid}}$ **(Page-oriented nested-loop join)**

**250 Pages**

**10 Pages**

**(Scan; Write to temp T1)** $\sigma_{\text{gpa}>3.0}$

$\sigma_{\text{cid}=3220}$ **(Scan; Write to temp T2)**

**Students**

**CourseEnrolled**

Cost of $\sigma_{\text{cid}=3220}$ (E)= Cost to scan E + Cost to write T2
= 500 pages + size(T2)
Since there are 50 courses (i.e., 50 cids) in the Table E and these cids are uniformly distributed in this table (Assumption 2), we have:
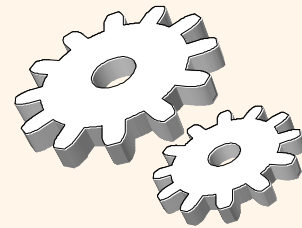size(T2) = 500/50 = 10 pages

# *Cost of Evaluating the Alternative Plan 1*

- ❖ Cost of *page-oriented nested loop join* of T1 and T2
  - Cost = $10 + 10 \times 250 = 2510$ I/Os

- ❖ Total cost of Alternative Plan 1

  = cost of selection operations + cost of join operation

  = (1250 + 510) + 2510

  = 1760 + 2510

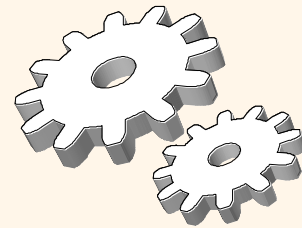  = 4270 I/Os     **(better than that of the original QEP)**

# *Question 2*

❖ Consider the relations A(<u>a</u>, b, c) and B(<u>a</u>, x, y) and the following SQL query.

SELECT *
FROM A, B
WHERE b > 20 AND A.a = B.a

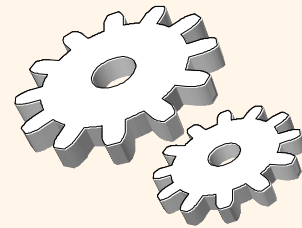a) Write down the QEP for this SQL query.
b) Write down another QEP for this SQL query.

# *Question 3*

Given a relational algebra expression:

- $\pi_{\text{sid, sname}}(S) \bowtie_{\text{S.sid=E.sid}} \sigma_{\text{cid=3220}}(E)$

1) What is its equivalent relational algebra expression?
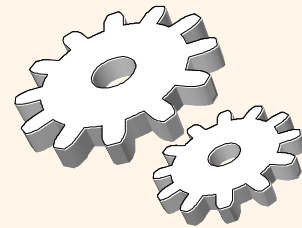2) Draw the QEPs for these two relational algebra expressions.

# Question 4

❖ Given three relations Students(SID, Name, Email, GPA), Courses(CID, Name, Day), and EnrolledCourses(SID, CID, Date) *underline: without indexes*, and the following SQL query,

SELECT  Students.Name, Students.Email
FROM Students, Courses, EnrolledCourses
WHERE Students.SID=EnrolledCourses.SID
        AND Courses.CID=EnrolledCourses.CID
        AND Courses.CID=7640
        AND Students.GPA>3.4

draw three possible query evaluation plans for solving this query. (Use S, C, and E to denote relations Students, Courses, and EnrolledCourses, respectively)

# *Solution to Question 1*

❖ Using distributive rule 2 for selection & join

$$\sigma_{\theta_1 \wedge \theta_2}(R \bowtie_\theta S) \Leftrightarrow \sigma_{\theta_1}(R) \bowtie_\theta \sigma_{\theta_2}(S) \ , \ \theta_1 \subseteq S, \ \theta_2 \subseteq R$$

$$\pi_{\text{sname,gpa,cid}} \left( \sigma_{\text{gpa}>3.0 \ \wedge \ \text{cid}=3220}(S \bowtie_{\text{S.sid=E.sid}} E) \right)$$

$$\Downarrow$$

$$\pi_{\text{sname,gpa,cid}} \left( \sigma_{\text{gpa}>3.0}(S) \bowtie_{\text{S.sid=E.sid}} \sigma_{\text{cid}>3220}(E) \right)$$

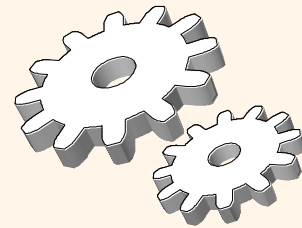❖ Using commutative rule for selection & projection

$$\pi_L(\sigma_\theta(R)) \Leftrightarrow \sigma_\theta(\pi_L(R)), \ \theta \subseteq L$$

$$\pi_{\text{sname,gpa,cid}} \left( \sigma_{\text{gpa}>3.0 \ \wedge \ \text{cid}=3220}(S \bowtie_{\text{S.sid=E.sid}} E) \right)$$

$$\Downarrow$$

$$\sigma_{\text{gpa}>3.0 \ \wedge \ \text{cid}=3220} \left( \pi_{\text{sname,gpa,cid}}(S \bowtie_{\text{S.sid=E.sid}} E) \right)$$

# *Solution to Question 2*

a) $\sigma_{b>20}(A \bowtie_{A.a=B.a} B)$     b) $\sigma_{b>20}(A) \bowtie_{A.a=B.a} B$

$\sigma_{b>20}$ **(Scan)**

$\bowtie_{A.a=B.a}$ **(Page-oriented nested-loop join)**

**A**     **B**

$\bowtie_{T1.a=B.a}$ **(Page-oriented nested-loop join)**

$\sigma_{b>20}$ **(Scan; write to temp T1)**

**B**

**A**

37

# *Solution to Question 3*

❖ By Distributive Rule 1

$\sigma_{\text{cid}=3220}(\pi_{\text{sid, sname}}(S) \bowtie_{\text{S.sid=E.sid}} E)$

# *Solution to Question 3*

❖ $\pi_{\text{sid, sname}}(S) \bowtie_{\text{S.sid=E.sid}} \sigma_{\text{cid=3220}}(E)$

$\bowtie_{\text{S.sid=C.sid}}$ **(Page-oriented nested-loop join)**

**(Scan; write to temp T1)** $\pi_{\text{sid, sname}}$      $\sigma_{\text{cid=3220}}$ **(Scan; write to temp T2)**

**S**                    **C**

# Solution to Question 3

❖ $\sigma_{\text{CID}=3220}(\pi_{\text{SID, Name}}(S) \bowtie_{\text{S.SID=E.SID}} E)$

$\sigma_{\text{cid}=3220}$ **(On the fly)**

$\bowtie_{\text{S.sid=C.sid}}$ **(Page-oriented nested-loop join)**

**(Scan; write to temp T1)** $\pi_{\text{sid, sname}}$

**S**

**C**

# *Solution to Question 4*

❖ Relational algebra expression

$$\pi_{\text{S.Name, S.Email}}\left(\sigma_{\text{S.GPA>3.4} \land \text{C.CID=7640}}\left(S \bowtie_{\text{S.SID=E.SID}} \left(C \bowtie_{\text{C.CID=E.CID}} E\right)\right)\right)$$

**Or** $\pi_{\text{S.Name, S.Email}}\left(\sigma_{\text{S.GPA>3.4} \land \text{C.CID=7640}}\left(\left(S \bowtie_{\text{S.SID=E.SID}} E\right) \bowtie_{\text{C.CID=E.CID}} C\right)\right)$

❖ Enumerate alternative relational algebra expressions based on the 1$^{\text{st}}$ one

**Distributive Rule 1**

$$\pi_{\text{S.Name, S.Email}}\left(\sigma_{\text{C.CID=7640}}\left(\sigma_{\text{S.GPA>3.4}}(S) \bowtie_{\text{S.SID=E.SID}} \left(C \bowtie_{\text{C.CID=E.CID}} E\right)\right)\right)$$

$$\pi_{\text{S.Name, S.Email}}\left(\sigma_{\text{S.GPA>3.4}}\left(S \bowtie_{\text{S.SID=E.SID}} \sigma_{\text{C.CID=7640}}\left(C \bowtie_{\text{C.CID=E.CID}} E\right)\right)\right)$$

**Can further apply Distributive Rules**

**Distributive Rule 2**

$$\pi_{\text{S.Name, S.Email}}\left(\sigma_{\text{S.GPA>3.4}}(S) \bowtie_{\text{S.SID=E.SID}} \sigma_{\text{C.CID=7640}}\left(C \bowtie_{\text{C.CID=E.CID}} E\right)\right)$$

**Can further apply Distributive Rules**

# *Solution to Question 4*

❖ QEP 1

$$\pi_{\text{S.Name, S.Email}}\left(\sigma_{\text{S.GPA}>3.4 \,\wedge\, \text{C.CID}=7640}\left(S \bowtie_{\text{S.SID}=\text{E.SID}} \left(C \bowtie_{\text{C.CID}=\text{E.CID}} E\right)\right)\right)$$

$\pi_{\text{S.Name, S.Email}}$ **(On the fly)**

$\sigma_{\text{S.GPA}>3.4 \wedge \text{C.CID}=7640}$ **(On the fly)**

$\bowtie_{\text{S.SID}=\text{C.SID}}$ **(Page-oriented nested-loop join)**

**S**

$\bowtie_{\text{C.CID}=\text{E.CID}}$ **(Page-oriented nested-loop join)**

**C**   **E**

# *Solution to Question 4*

❖ QEP 2

$$\pi_{\text{S.Name, S.Email}}\left(\sigma_{\text{C.CID}=7640}(\sigma_{\text{S.GPA}>3.4}\ (S)\ \bowtie_{\text{S.SID}=\text{E.SID}}\ (C\ \bowtie_{\text{C.CID}=\text{E.CID}}\ E))\right)$$

$\pi_{\text{S.Name, S.Email}}$  **(On the fly)**

$\sigma_{\text{C.CID}=7640}$  **(On the fly)**

$\bowtie_{\text{S.SID}=\text{C.CID}}$  **(Page-oriented nested-loop join)**

**(Scan; write to temp T1)**  $\sigma_{\text{S.GPA}>3.4}$  $\bowtie_{\text{C.CID}=\text{E.CID}}$  **(Page-oriented nested-loop join)**

**S**  **C**  **E**

# *Solution to Question 4*

❖ QEP 3

$$\pi_{\text{S.Name, S.Email}}(\sigma_{\text{S.GPA>3.4}}(S \bowtie_{\text{S.SID=E.SID}} \sigma_{\text{C.CID=7640}}(C \bowtie_{\text{C.CID=E.CID}} E)))$$

$\pi_{\text{S.Name, S.Email}}$    **(On the fly)**

$\sigma_{\text{S.GPA>3.4}}$    **(On the fly)**

$\bowtie_{\text{S.SID=C.SID}}$    **(Page-oriented nested-loop join)**

**S**

$\sigma_{\text{C.CID=7640}}$    **(Scan; write to temp T1)**

$\bowtie_{\text{C.CID=E.CID}}$    **(Page-oriented nested-loop join)**

**C**   **E**

# *Solution to Question 4*

❖ QEP 4

$$\pi_{\text{S.Name, S.Email}}(\sigma_{\text{S.GPA>3.4}}(S) \bowtie_{\text{S.SID=E.SID}} \sigma_{\text{C.CID=7640}}(C \bowtie_{\text{C.CID=E.CID}} E))$$

$\pi_{\text{S.Name, S.Email}}$   **(On the fly)**

$\bowtie_{\text{S.SID=C.SID}}$   **(Page-oriented nested-loop join)**

**(Scan; write to temp T1)** $\sigma_{\text{S.GPA>3.4}}$     $\sigma_{\text{C.CID=7640}}$ **(Scan; write to temp T1)**

**S**

$\bowtie_{\text{C.CID=E.CID}}$ **(Page-oriented nested-loop join)**

**C**   **E**