

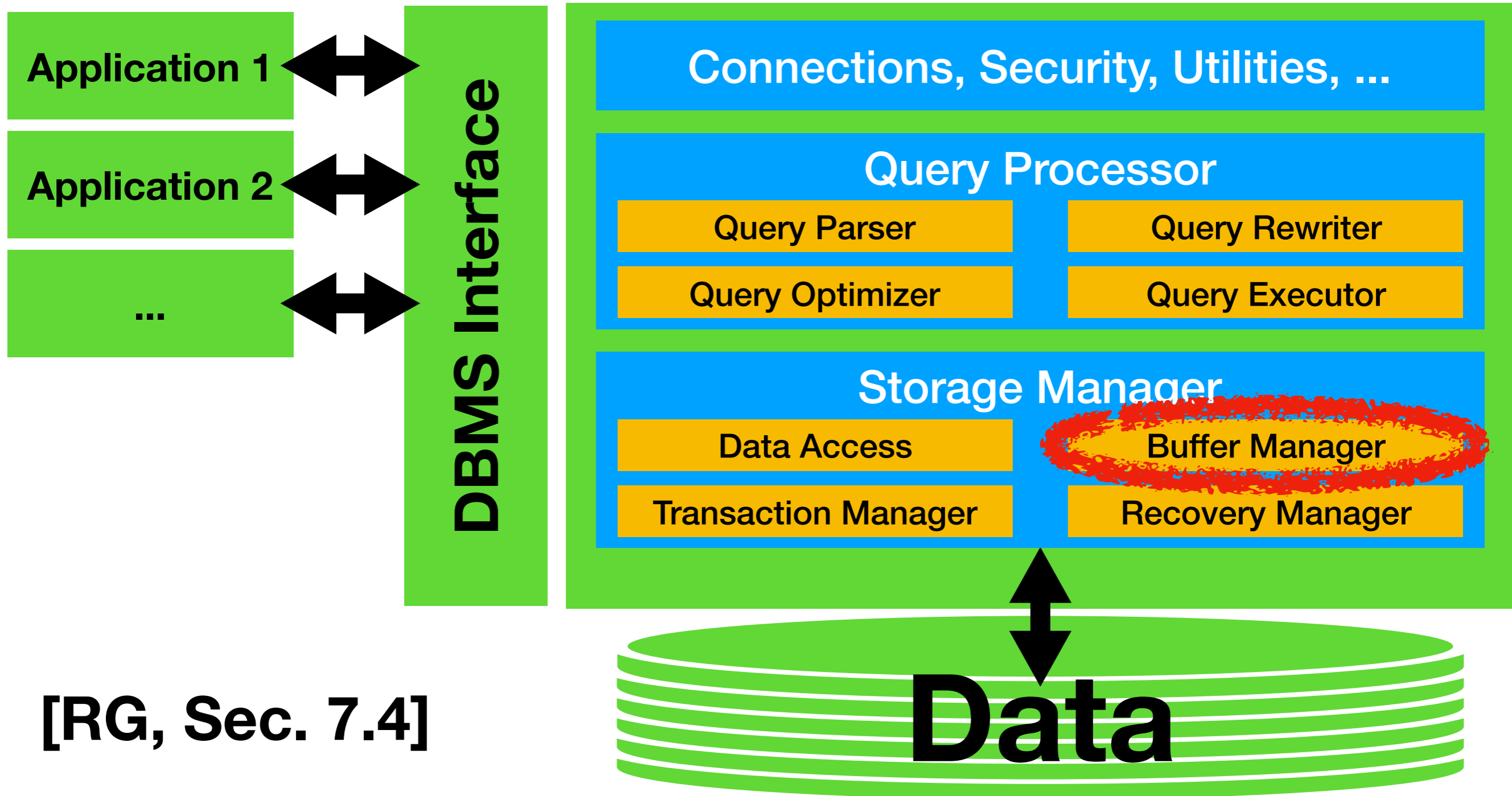
# Query Processing Overview

Immanuel Trummer

[itrummer@cornell.edu](mailto:itrummer@cornell.edu)

[www.itrummer.org](http://www.itrummer.org)

# Database Management Systems (DBMS)



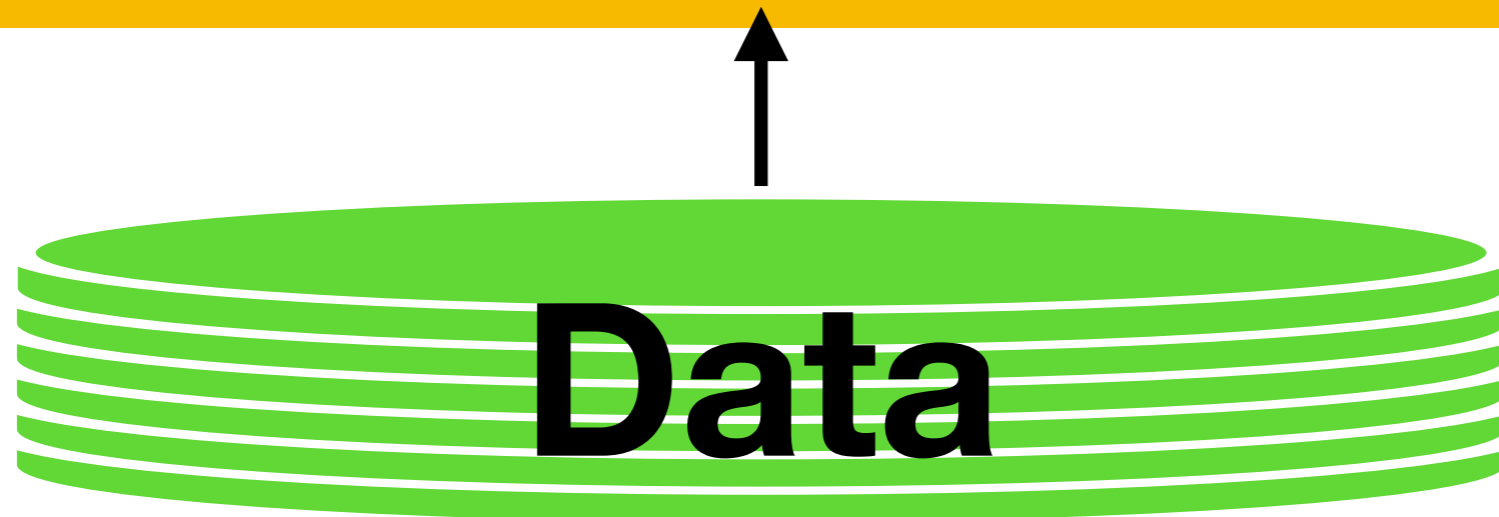
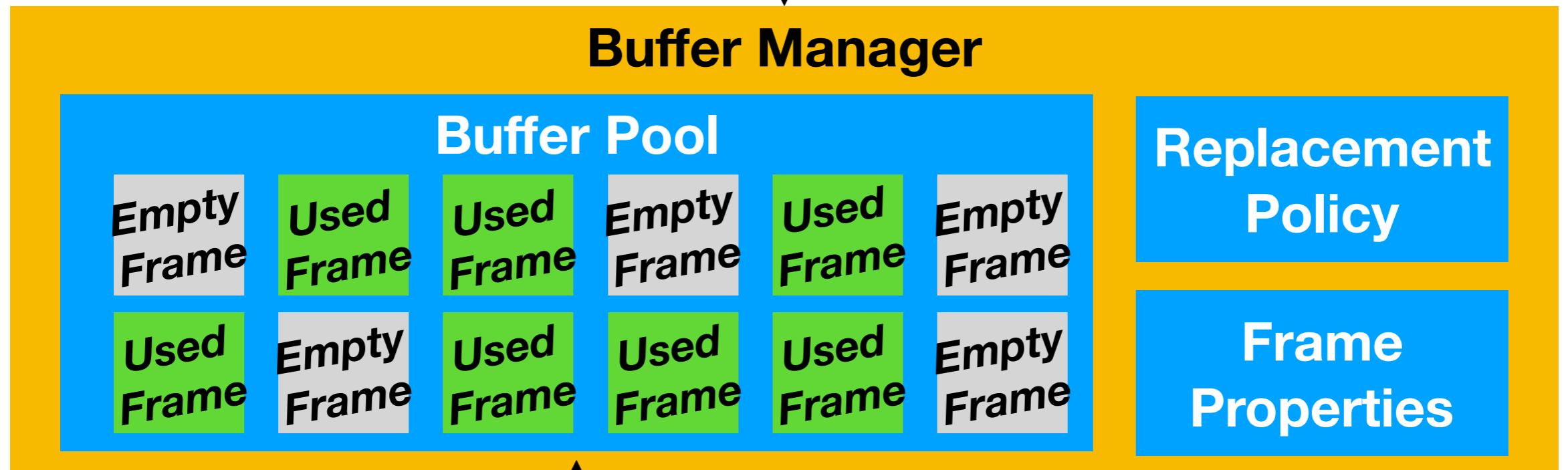
[RG, Sec. 7.4]

# Buffer Manager

- Decides when to **move data** between disk and RAM
- Tries to **reduce** data movements using heuristics
- Buffer manager manages "**buffer pool**"
  - Buffer pool: main memory **reserved** for DBMS
  - Divided into page-sized slots called "**frames**"
  - Stores meta-data about each slot

# Buffer Manager Illustration

*Requests for specific pages*



# Frame Properties

- **Page ID**: which page is currently stored in frame?
  - Allows **matching** page requests to frames
- **Pin count**: how many processes are using a page?
  - Can only **evict** page if pin count reaches zero
- **Dirty bit**: in-memory page deviates from disk version?
  - Must **write** page to disk before evicting it

# Processing Page Requests

- Case 1: **Cache Hit** (requested page cached)
  - **Increase pin** count and **return page** address
- Case 2: **Cache Miss** (requested page not cached)
  - **Choose frame** for replacement (replacement policy)
  - If frame contains dirty page then **write** it to disk
  - **Read** requested page from disk and store in frame
  - **Increase pin** count and **return page** address

# Why Not Rely on OS?

- **OS caches** pages as well (virtual memory)
- Why do we want a **separate** buffer manager?
  - DBMS **knows** its access patterns ahead of time
    - Can exploit for smarter **replacements**
  - DBMS must control **page writes** for safety guarantees

# LRU Replacement Policy

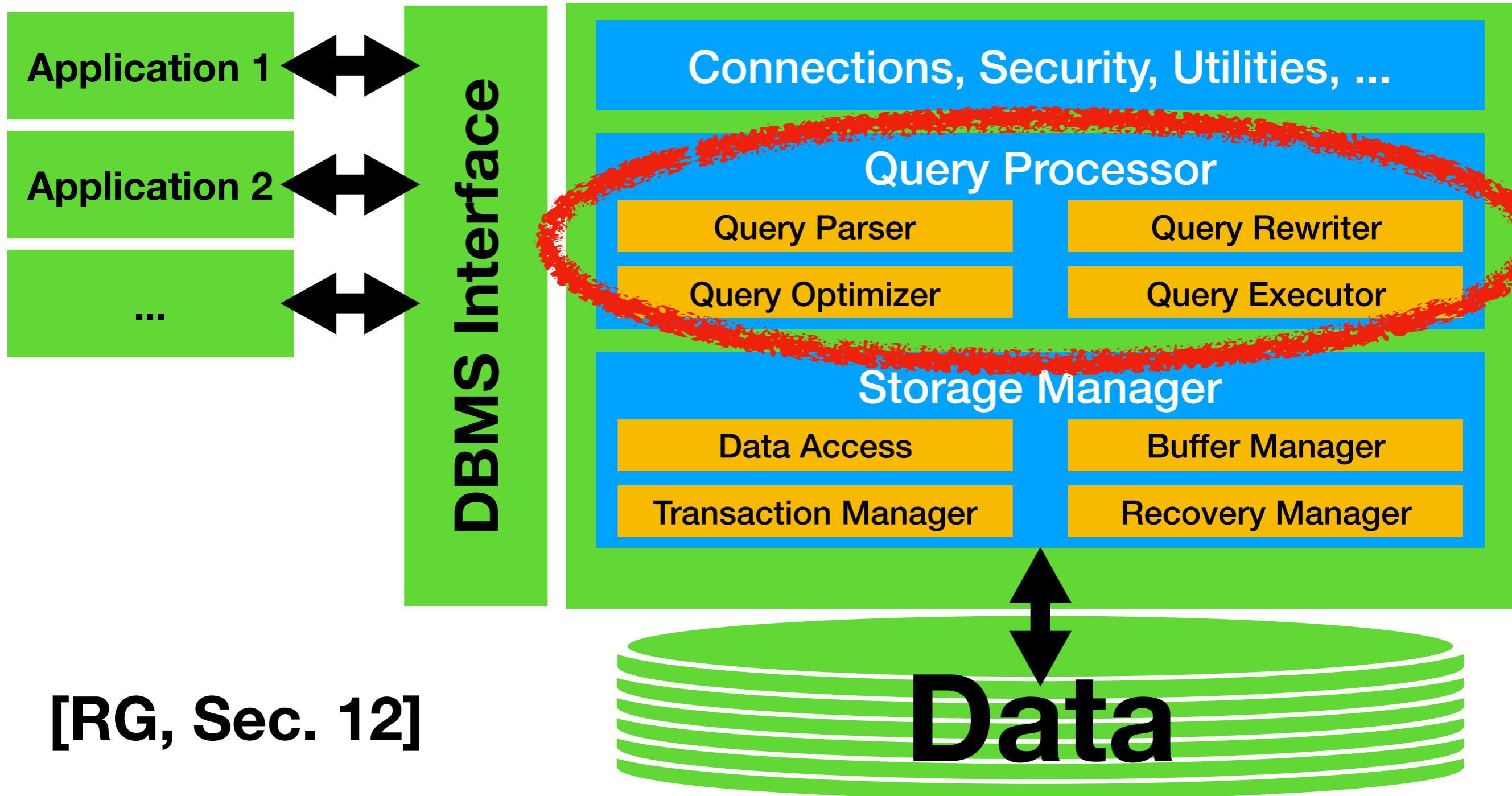
- Want to replace page required **farthest** in the future
  - Doing so **reduces** expensive cache misses
- However: **difficult** to predict that in general
- Heuristic: remove **least recently used** page (LRU)
  - Did not use page for long time, **unlikely to do soon**



# Sequential Flooding

- DBMS often have **particular** access patterns
- For instance: **keep scanning** pages in round robin mode
- Least recently page is used again **soonest**
- Makes **LRU** policy highly **sub-optimal**
- Otherwise a reasonable strategy!

# Database Management Systems (DBMS)



[RG, Sec. 12]

# Query Processing

- Input query is parsed (**Parser**) and simplified (**Rewriter**)
- **Query optimizer** generates optimized execution plan
- Executing plan (**Executor**) produces query result

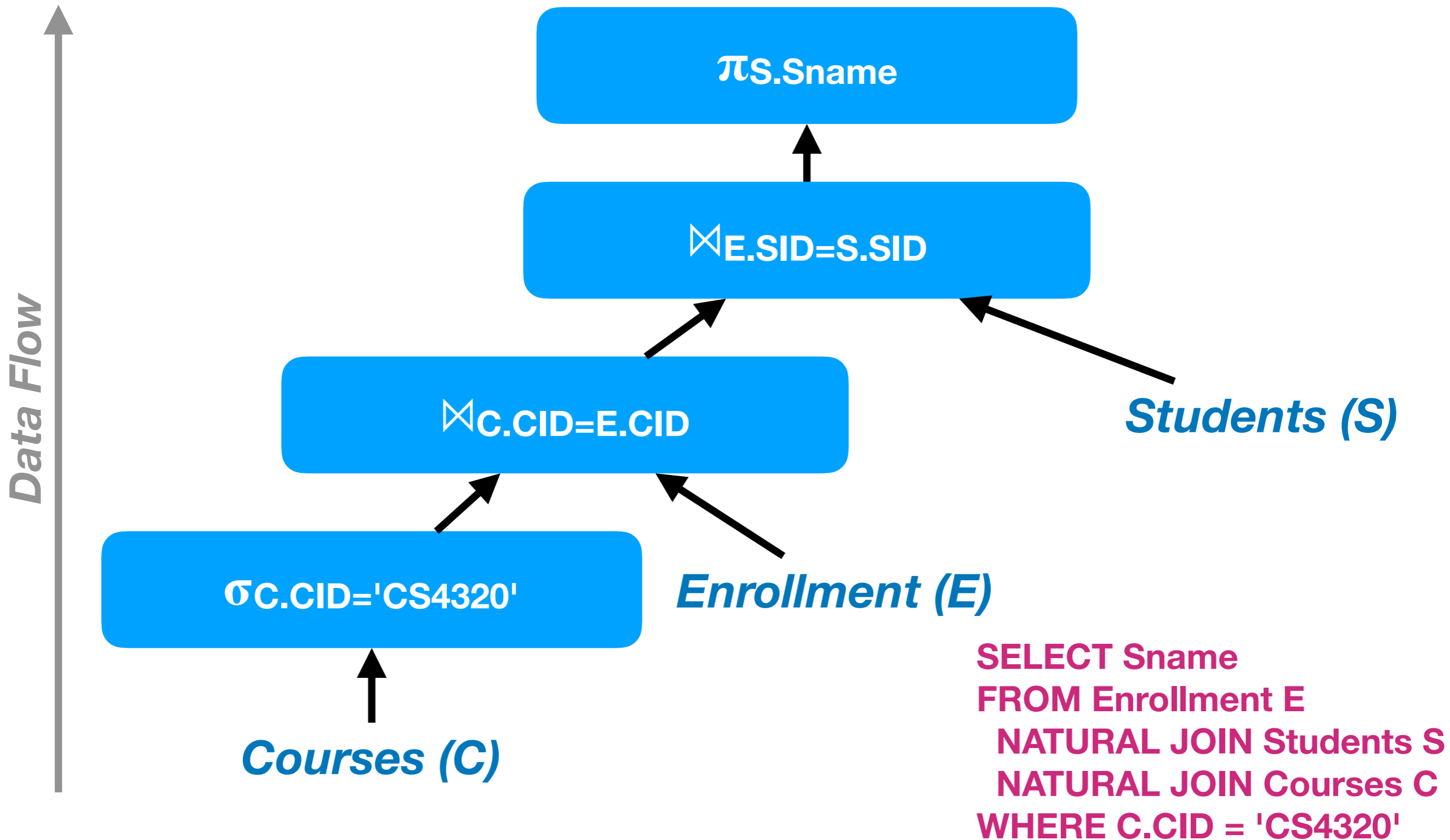
# Query Plans

- Describe **how** to generate required data
- Typically represented as a **tree**
- Each leaf node represents a database **table**
- Each inner node represents an **operation**
- Tree edges represent **data flow**

# Operators

- Query plans use fixed set of **standard operators**
- Consumes relation(s) and produces one **relation**
- **Filter operator ( $\sigma$ )**: discard rows based on condition
- **Projection operator ( $\pi$ )**: discard columns
- **Join operator ( $\bowtie$ )**: find matching tuple pairs
- ...

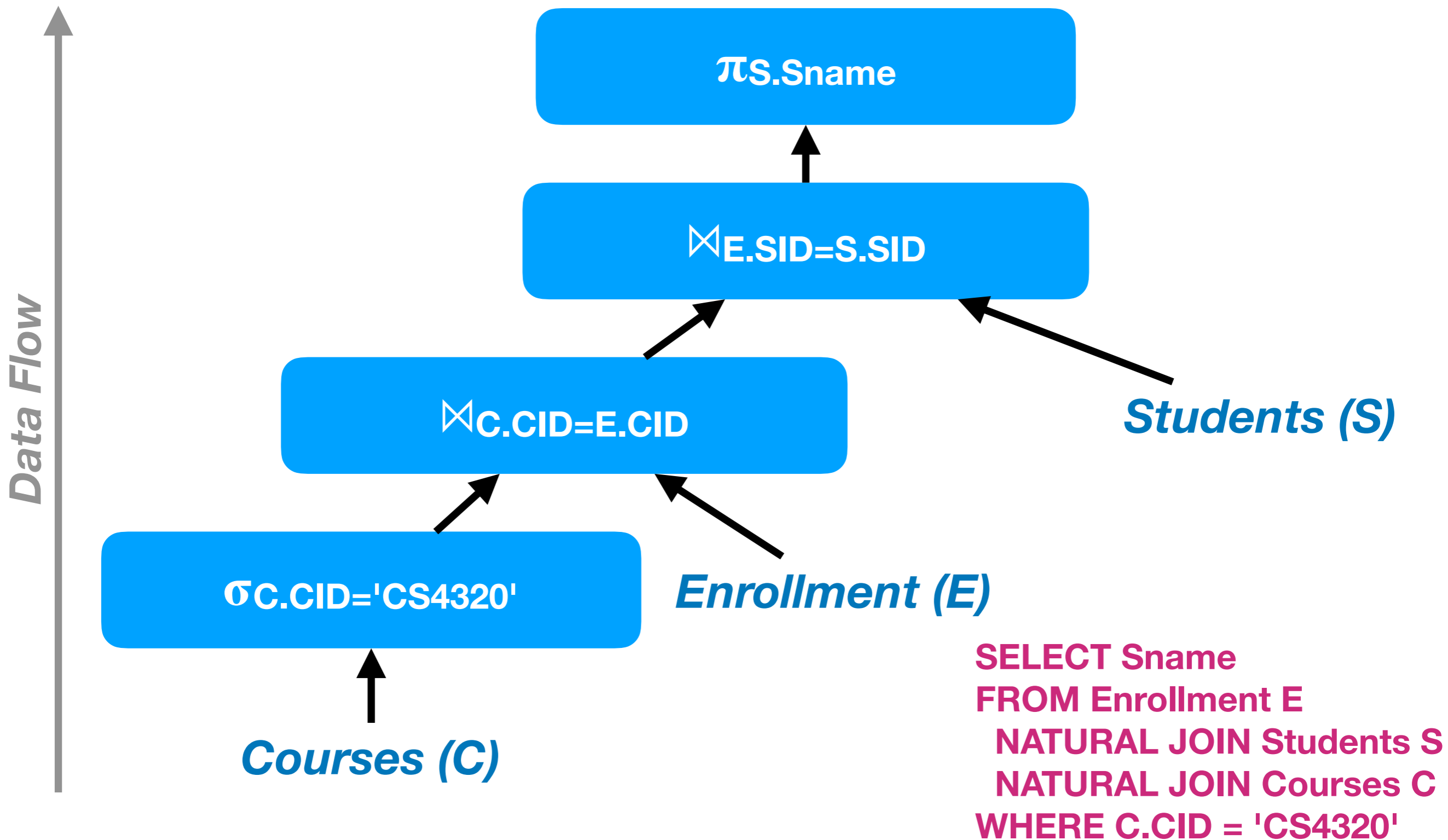
# Example Plan



# Relational Algebra

- Mathematical **foundations** for describing query plans
- Represents query plan as "mathematical **expression**"
- Won't discuss in too much detail here

# Example Plan



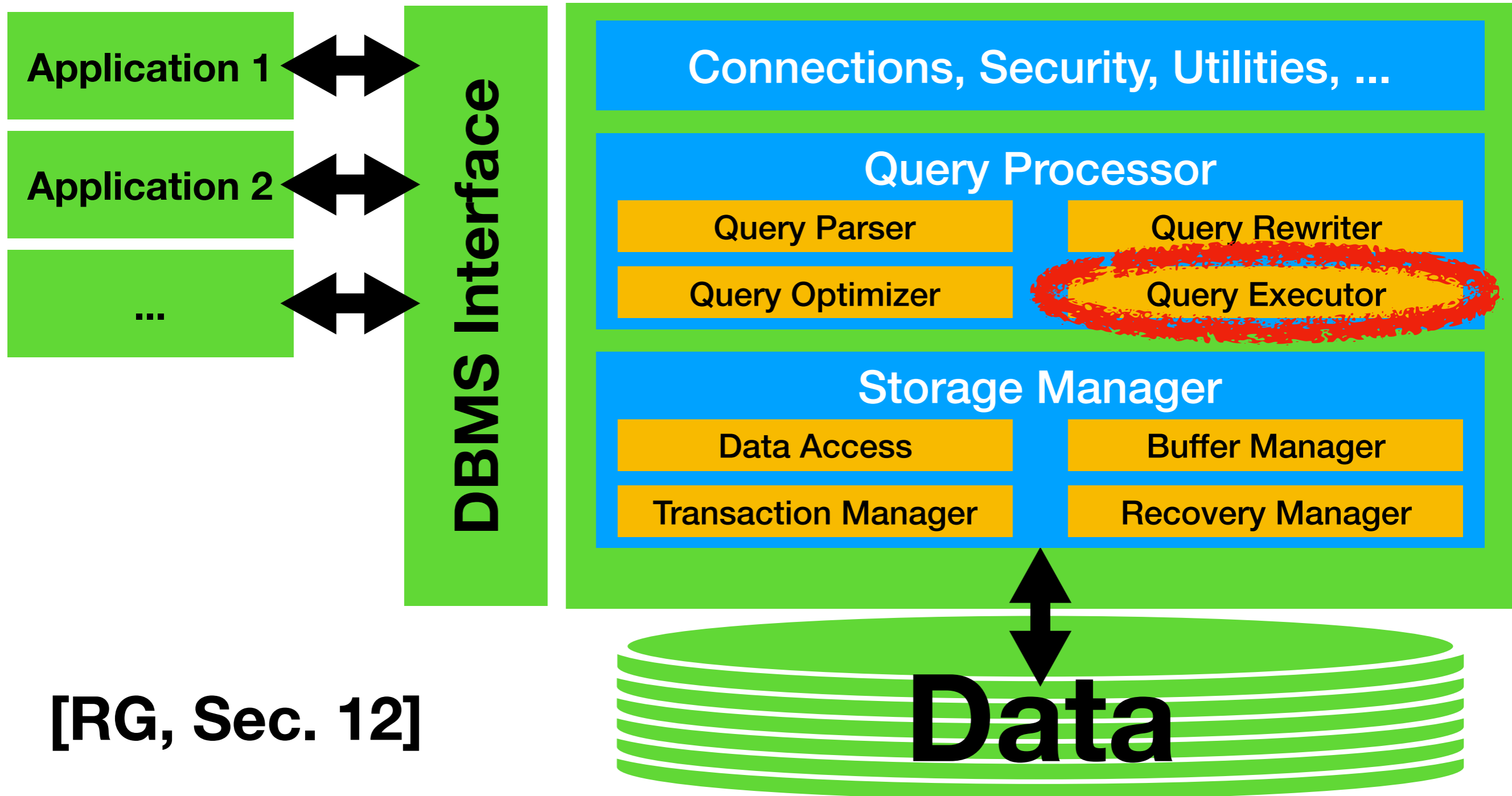


# Relational Algebra Representation

$\pi_{Sname}(\sigma_{C.CID='CS4320'}(C) \bowtie_{C.CID=E.CID} E \bowtie_{E.SID=S.SID} S)$

*(Evaluated from inner to outer expressions)*

# Database Management Systems (DBMS)



[RG, Sec. 12]

# Outlook

- Will discuss how to implement **standard operators**
- Often have **multiple** implementations of same operator
- Can choose **most efficient** implementation at each point
- Will also discuss about **cost estimation**
  - Assumption: cost  $\sim$  **number of pages** read/written

# Filter Operator ( $\sigma$ )

# How to Filter?

- Want to retrieve table rows satisfying a **predicate**
- Simplest option: **scan** all pages, check each entry
- Option if sorted: **binary search** for specific predicates
- Can use **indexes** if available (right table, right key)

# Costing Example

**SELECT \* FROM Enrollment WHERE CID='CS4320'**

Property	Value
Nr. Students	60,000
Nr. Enrollments/Student	10
Size/Enrollment Entry	10 Bytes
Bytes/Page	1,000 Bytes

***Calculate scan cost!***

Nr. Entries/Page	?
Nr. Enrollment Pages	?
Total Scan Cost	?

# Costing Example

**SELECT \* FROM Enrollment WHERE CID='CS4320'**

Property	Value
Nr. Students	60,000
Nr. Enrollments/Student	10
Size/Enrollment Entry	10 Bytes
Bytes/Page	1,000 Bytes

***Calculate scan cost!***

Nr. Entries/Page	100
Nr. Enrollment Pages	6,000
Total Scan Cost	6,000

*What About  
Output Cost?*



# Costing Example

**SELECT \* FROM Enrollment WHERE CID='CS4320'**

Property	Value
Nr. Students	60,000
Nr. Courses	100
Nr. Enrollments/Student	10
Size/Enrollment Entry	10 Bytes
Bytes/Page	1,000 Bytes

***Sorted by CID - Calculate binary search cost!***

Nr. Entries/Page	100
Nr. Enrollment Pages	6,000
Nr. Search Steps	?
Nr. Pages to Scan	?
Total Cost	?

# Costing Example

**SELECT \* FROM Enrollment WHERE CID='CS4320'**

Property	Value
Nr. Students	60,000
Nr. Courses	100
Nr. Enrollments/Student	10
Size/Enrollment Entry	10 Bytes
Bytes/Page	1,000 Bytes

***Sorted by CID - Calculate binary search cost!***

Nr. Entries/Page	100
Nr. Enrollment Pages	6,000
Nr. Search Steps	13
Nr. Pages to Scan	60
Total Cost	~ 73

*Where Did We  
Simplify?*

# Costing Example

**SELECT \* FROM Enrollment WHERE CID='CS4320'**

Property	Value
Nr. Students	60,000
Nr. Courses	100
Nr. Enrollments/Student	10
Size/Enrollment Entry	10 Bytes
Bytes/Page	1,000 Bytes
Index Fanout	100

***Tree Index with Data on CID - Calculate Access Cost!***

Nr. Entries/Page	100
Nr. Enrollment Pages	6,000
Nr. Inner Node Visits	?
Nr. Leaf Node Visits	?
Total Cost	?

# Costing Example

**SELECT \* FROM Enrollment WHERE CID='CS4320'**

Property	Value
Nr. Students	60,000
Nr. Courses	100
Nr. Enrollments/Student	10
Size/Enrollment Entry	10 Bytes
Bytes/Page	1,000 Bytes
Index Fanout	100

***Tree Index with Data on CID - Calculate Access Cost!***

Nr. Entries/Page	100
Nr. Enrollment Pages	6,000
Nr. Inner Node Visits	2
Nr. Leaf Node Visits	60
Total Cost	62

# Costing Example

**SELECT \* FROM Enrollment WHERE CID='CS4320'**

Property	Value
Nr. Students	60,000
Nr. Courses	100
Nr. Enrollments/Student	10
Size/Enrollment Entry	10 Bytes
Bytes/Page	1,000 Bytes
Index Fanout	100

***Unclustered Tree Index on CID - Calculate Access Cost!***

Nr. Entries/Page	100
Nr. Enrollment Pages	6,000
Nr. Inner Node Visits	2
Nr. Leaf Node Visits	60
Nr. Data Pages Read	6,000
Total Cost	6,062

# Example Summary

<b>Scan Cost</b>	6,000
<b>Binary Search</b>	73
<b>Index with Data</b>	62
<b>Unclustered Index</b>	6,062

# Insights

- Index or sort orders **can** speed up filtering
- However, may **not always** be more efficient
- Need to **calculate cost** of alternatives and compare
  - This is what the **query optimizer** does ...



# Multi-Predicate Filtering

- May have to retrieve entries satisfying **two predicates**
- **Scanning** all pages always works
- Can **use index** for first predicate, then **check** second
- Could **merge results** from two indices for both predicates