# LAST CLASS

Course objectives and expectations.
→ I will assign note taking schedule tonight.

Motivation for why query optimization is important and a hard problem.

# TODAY'S AGENDA

Background

Heuristics

Heuristics + Cost-based Search

# THE DAWN OF THE RELATIONAL MODEL

In the late 1960s, early DBMSs required developers to write queries using procedural code.
→ Example: CODASYL

The developer had to choose access paths and execution ordering based on the current database contents.
→ If the database changes, then the developer must rewrite the query code.

# THE DAWN OF THE RELATIONAL MODEL

In the late 1960s, early DBMSs required developers to write queries using procedural code.
→ Example: CODASYL

The developer had to choose access paths and execution ordering based on the current database contents.
→ If the database changes, then the developer must rewrite the query code.

In order to focus the role of programmer as navigator, let us enumerate his opportunities for record access. These represent the commands that he can give to the database system—singly, multiply or in combination with each other-- as he picks his way through the data to resolve an inquiry or to complete an update.

1. He can start at the beginning of the database, or at any known record, and sequentially access the "next" record in the database until he reaches a record of interest or reaches the end.

2. He can enter the database with a database key that provides direct access to the physical location of a record. (A database key is the permanent virtual memory address assigned to a record at the time that it was created.)

3. He can enter the database in accordance with the value of a primary data key. (Either the indexed sequential or randomized access techniques will yield the same result.)

4. He can enter the database with a secondary data key value and sequentially access all records having that particular data value for the field.

5. He can start from the owner of a set and sequentially access all the member records. (This is equivalent to converting a primary data key into a secondary data key.)

6. He can start with any member record of a set and access either the next or prior member of that set.

7. He can start from any member of a set and access the owner of the set, thus converting a secondary data key into a primary data key.

# THE DAWN OF THE RELATIONAL MODEL

In the late 1960s, early DBMSs required developers to write queries using procedural code

→ Example: CODASYL

The developer had to ~~specify the~~ paths and execution ~~order based~~ on the current database contents.

→ If the database changes, then the developer must rewrite the query code.

In order to focus the role of programmer as navigator, let us enumerate his opportunities for record access. These represent the commands that he can give to the database system—singly, multiply or in combination with each other-- as he picks his way through the data to resolve an inquiry or to complete an update.

1. He can start at the beginning of the database, or at any known record, and sequentially access the "next" record in the ~~database~~

Each of these access methods is interesting in itself, and all are very useful. However, ==it is the synergistic usage of the entire collection which gives the programmer great and expanded powers to come and go within a large database while accessing only those records of interest in responding to inquiries and updating the database in anticipation of future inquiries.==

~~...~~ and sequentially access all records having that particular data value for the field.

5. He can start from the owner of a set and sequentially access all the member records. (This is equivalent to converting a primary data key into a secondary data key.)

6. He can start with any member record of a set and access either the next or prior member of that set.

7. He can start from any member of a set and access the owner of the set, thus converting a secondary data key into a primary data key.

# THE DAWN OF THE RELATIONAL MODEL

In the late 1960s, early DBMSs required developers to write queries using procedural code.
→ Example: CODASYL

The developer had to choose access paths and execution ordering based on the current database contents.
→ If the database changes, then the developer must rewrite the query code.

*Retrieve the names of artists that appear on the DJ Mooshoo Tribute mixtape.*

```
PROCEDURE GET_ARTISTS_FOR_ALBUM;
BEGIN
    /* Declare variables */
    DECLARE ARTIST_RECORD ARTIST;
    DECLARE APPEARS_RECORD APPEARS;
    DECLARE ALBUM_RECORD ALBUM;

    /* Start navigation */
    FIND ALBUM USING ALBUM.NAME = "Mooshoo Tribute"
        ON ERROR DISPLAY "Album not found" AND EXIT;

    /* For each appearance on the album */
    FIND FIRST APPEARS WITHIN APPEARS_ALBUM OF ALBUM_RECORD
        ON ERROR DISPLAY "No artists found for this album" AND EXIT;

    /* Loop through the set of APPEARS */
    REPEAT
        /* Navigate to the corresponding artist */
        FIND OWNER WITHIN ARTIST_APPEARS OF APPEARS_RECORD
            ON ERROR DISPLAY "Error finding artist";
        /* Display artist name */
        DISPLAY ARTIST_RECORD.NAME;
        /* Move to the next APPEARS record in the set */
        FIND NEXT APPEARS WITHIN APPEARS_ALBUM OF ALBUM_RECORD
            ON ERROR EXIT;
    END REPEAT;
END PROCEDURE;
```

# THE DAWN OF THE RELATIONAL MODEL

In the late 1960s, early DBMSs required developers to write queries using procedural code.
→ Example: CODASYL

The developer had to choose access paths and execution ordering based on the current database contents.
→ If the database changes, then the developer must rewrite the query code.

*Retrieve the names of artists that appear on the DJ Mooshoo Tribute mixtape.*

```
PROCEDURE GET_ARTISTS_FOR_ALBUM;
BEGIN
  /* Declare variables */
  DECLARE      ST_RECORD ARTIST;
  DECLA         S_RECORD APPEARS;
  DECLA         CORD ALBUM;

  /* Start        gation
  FIND ALBUM USING
      ON ERROR DI

  /* For each a
  FIND FIRST A      WITHIN A    BUM OF AL    RECORD
      ON ERROR      LAY "No arti    ound for thi    um" AND EXIT;

  /* Loop thro     he set of A      */
  REPEAT
      /* Naviga            ond
      FIND OWNER             _APPEAR            RECORD
          ON ERROR                        ";
      /* Dis    artis
      DISP       ST_RECOR
      /*      e next        in the
      F       APPEARS WI    N A    ARS    BUM OF AL    D
               RROR EXIT;
  END REPEAT;
END PROCEDURE;
```

# THE DAWN OF THE RELATIONAL MODEL

In the late 1960s, early DBMSs required developers to write queries using procedural code.
→ Example: <u>CODASYL</u>

The developer had to choose access paths and execution ordering based on the current database contents.
→ If the database changes, then the developer must rewrite the query code.

*Retrieve the names of artists that appear on the DJ Mooshoo Tribute mixtape.*

```
SELECT ARTIST.NAME
  FROM ARTIST, APPEARS, ALBUM
 WHERE ARTIST.ID=APPEARS.ARTIST_ID
   AND APPEARS.ALBUM_ID=ALBUM.ID
   AND ALBUM.NAME="Mooshoo Tribute"
```

# RELATIONAL MODEL

**Structure:** The definition of the database's relations and their contents independent of their physical representation.

**Integrity:** Ensure the database's contents satisfy constraints.

**Manipulation:** Declarative API for accessing and modifying a database's contents via sets.

# RELATIONAL MODEL

Early relational DBMS implementations:
→ **Peterlee Relational Test Vehicle** – IBM Research (UK)
→ **System R** – IBM Research (San Jose)
→ **INGRES** – U.C. Berkeley
→ **Oracle** – Larry Ellison
→ **Mimer** – Uppsala University



*Gray*



*Stonebraker*



*Ellison*

# HISTORY OF QUERY OPTIMIZERS

**Choice #1: Heuristics**
→ INGRES (1970s), Oracle (until mid 1990s)

**Choice #2: Heuristics + Cost-based Join Search**
→ System R (1970s), early IBM DB2

**Choice #3: Stratified Search**
→ IBM STARBURST (late 1980s), now IBM DB2 + Oracle

**Choice #4: Unified Search**
→ Volcano/Cascades in 1990s, now MSSQL + Greenplum

**Choice #5: Randomized Search**
→ Academics in the 1980s, current Postgres

# HEURISTIC-BASED OPTIMIZATION

Define static rules that transform logical operators to a physical plan <u>without</u> a cost model.
→ Perform most restrictive selection early
→ Perform all selections before joins
→ Predicate/Limit/Projection pushdowns
→ Join ordering based on simple rules or cardinality estimates

*Stonebraker*

**Examples**: INGRES (until mid-1980s) and Oracle (until early-1990s), MongoDB, most new DBMSs.

QUERY PROCESSING IN A RELATIONAL DATABASE
MANAGEMENT SYSTEM
*VLDB 1979*

# RELATIONAL ALGEBRA EQUIVALENCES

Two relational algebra expressions are <u>equivalent</u> if they generate the same set of tuples.

These equivalences allow the DBMS to manipulate and transform a query plan into different forms without effecting the correctness of its output.
→ This is how a heuristic-based optimizer identifies better query plans without a cost model.

# RELATIONAL ALGEBRA EQUIVALENCES

**Selections:**

→ Perform filters as early as possible.

→ Breakup a complex predicate into conjunctive clauses and push down to lowest part of plan as possible.

$$\sigma_{p1 \wedge p2 \wedge \dots pn}(R) = \sigma_{p1}(\sigma_{p2}(\dots \sigma_{pn}(R)))$$

Simplify complex predicates:

→ (X=Y AND Y=3) ➜ X=3 AND Y=3

→ (X=1+1) ➜ X=2

→ (X=YEAR('1/15/2025') ➜ X=2025

# RELATIONAL ALGEBRA EQUIVALENCES

**Joins:**

$\rightarrow$ Commutative:

$R \bowtie S = S \bowtie R$

$\rightarrow$ Associative:

$(R \bowtie S) \bowtie T = R \bowtie (S \bowtie T)$

The number of different join orderings for an n-way join is a **Catalan Number** ($\approx 4^n$)

$\rightarrow$ Exhaustive enumeration will be too slow.

# LOGICAL QUERY OPTIMIZATION

Split Conjunctive Predicates

Predicate Pushdown

Replace Cartesian Products with Joins

Projection Pushdown

# SPLIT CONJUNCTIVE PREDICATES

```
SELECT ARTIST.NAME
  FROM ARTIST, APPEARS, ALBUM
 WHERE ARTIST.ID=APPEARS.ARTIST_ID
   AND APPEARS.ALBUM_ID=ALBUM.ID
   AND ALBUM.NAME="Mooshoo Tribute"
```

Decompose predicates into their simplest forms to make it easier for the optimizer to move them around.



π ARTIST.NAME

σ ARTIST.ID=APPEARS.ARTIST_ID **AND**
APPEARS.ALBUM_ID=ALBUM.ID **AND**
ALBUM.NAME="Mooshoo Tribute"

ARTIST    APPEARS    ALBUM

# SPLIT CONJUNCTIVE PREDICATES

```
SELECT ARTIST.NAME
  FROM ARTIST, APPEARS, ALBUM
 WHERE ARTIST.ID=APPEARS.ARTIST_ID
   AND APPEARS.ALBUM_ID=ALBUM.ID
   AND ALBUM.NAME="Mooshoo Tribute"
```

Decompose predicates into their simplest forms to make it easier for the optimizer to move them around.



$\pi$ ARTIST.NAME

$\sigma$ ARTIST.ID=APPEARS.ARTIST_ID

$\sigma$ APPEARS.ALBUM_ID=ALBUM.ID

$\sigma$ ALBUM.NAME="Mooshoo Tribute"

ARTIST    APPEARS    ALBUM

# PREDICATE PUSHDOWN

```
SELECT ARTIST.NAME
  FROM ARTIST, APPEARS, ALBUM
 WHERE ARTIST.ID=APPEARS.ARTIST_ID
   AND APPEARS.ALBUM_ID=ALBUM.ID
   AND ALBUM.NAME="Mooshoo Tribute"
```

Move the predicate to the lowest point in the plan after Cartesian products.

# REPLACE CARTESIAN PRODUCTS

```
SELECT ARTIST.NAME
  FROM ARTIST, APPEARS, ALBUM
 WHERE ARTIST.ID=APPEARS.ARTIST_ID
   AND APPEARS.ALBUM_ID=ALBUM.ID
   AND ALBUM.NAME="Mooshoo Tribute"
```

Replace all Cartesian Products with inner joins using the join predicates.

# PROJECTION PUSHDOWN

```
SELECT ARTIST.NAME
  FROM ARTIST, APPEARS, ALBUM
 WHERE ARTIST.ID=APPEARS.ARTIST_ID
   AND APPEARS.ALBUM_ID=ALBUM.ID
   AND ALBUM.NAME="Mooshoo Tribute"
```

Eliminate redundant attributes before pipeline breakers to reduce materialization cost.

# INGRES OPTIMIZER

*Retrieve the names of people that appear on the DJ Mooshoo Tribute mixtape ordered by their artist id.*

```
SELECT ARTIST.NAME
  FROM ARTIST, APPEARS, ALBUM
 WHERE ARTIST.ID=APPEARS.ARTIST_ID
   AND APPEARS.ALBUM_ID=ALBUM.ID
   AND ALBUM.NAME="Mooshoo Tribute"
 ORDER BY ARTIST.ID
```

*Step #1: Decompose into single-value queries*

# INGRES OPTIMIZER

*Retrieve the names of people that appear on the DJ Mooshoo Tribute mixtape ordered by their artist id.*

```
SELECT ARTIST.NAME
  FROM ARTIST, APPEARS, ALBUM
 WHERE ARTIST.ID=APPEARS.ARTIST_ID
   AND APPEARS.ALBUM_ID=ALBUM.ID
   AND ALBUM.NAME="Mooshoo Tribute"
 ORDER BY ARTIST.ID
```

*Step #1: Decompose into single-value queries*

*Query #1*

```
SELECT ALBUM.ID AS ALBUM_ID INTO TEMP1
  FROM ALBUM
 WHERE ALBUM.NAME="Mooshoo Tribute"
```

*Query #2*

```
SELECT ARTIST.NAME
  FROM ARTIST, APPEARS, TEMP1
 WHERE ARTIST.ID=APPEARS.ARTIST_ID
   AND APPEARS.ALBUM_ID=TEMP1.ALBUM_ID
 ORDER BY APPEARS.ID
```

# INGRES OPTIMIZER

*Retrieve the names of people that appear on the DJ Mooshoo Tribute mixtape ordered by their artist id.*

```
SELECT ARTIST.NAME
  FROM ARTIST, APPEARS, ALBUM
 WHERE ARTIST.ID=APPEARS.ARTIST_ID
   AND APPEARS.ALBUM_ID=ALBUM.ID
   AND ALBUM.NAME="Mooshoo Tribute"
 ORDER BY ARTIST.ID
```

*Step #1: Decompose into single-value queries*

*Query #1*

```
SELECT ALBUM.ID AS ALBUM_ID INTO TEMP1
  FROM ALBUM
 WHERE ALBUM.NAME="Mooshoo Tribute"
```

*Query #2*

```
SELECT ARTIST.NAME
  FROM ARTIST, APPEARS, TEMP1
 WHERE ARTIST.ID=APPEARS.ARTIST_ID
   AND APPEARS.ALBUM_ID=TEMP1.ALBUM_ID
 ORDER BY APPEARS.ID
```

# INGRES OPTIMIZER

*Retrieve the names of people that appear on the DJ Mooshoo Tribute mixtape ordered by their artist id.*

```
SELECT ARTIST.NAME
  FROM ARTIST, APPEARS, ALBUM
 WHERE ARTIST.ID=APPEARS.ARTIST_ID
   AND APPEARS.ALBUM_ID=ALBUM.ID
   AND ALBUM.NAME="Mooshoo Tribute"
 ORDER BY ARTIST.ID
```

*Step #1: Decompose into single-value queries*

*Query #1*

```
SELECT ALBUM.ID AS ALBUM_ID INTO TEMP1
  FROM ALBUM
 WHERE ALBUM.NAME="Mooshoo Tribute"
```

*Query #3*

```
SELECT APPEARS.ARTIST_ID INTO TEMP2
  FROM APPEARS, TEMP1
 WHERE APPEARS.ALBUM_ID=TEMP1.ALBUM_ID
 ORDER BY APPEARS.ARTIST_ID
```

*Query #4*

```
SELECT ARTIST.NAME
  FROM ARTIST, TEMP2
 WHERE ARTIST.ARTIST_ID=TEMP2.ARTIST_ID
```

# INGRES OPTIMIZER

*Retrieve the names of people that appear on the DJ Mooshoo Tribute mixtape ordered by their artist id.*

```
SELECT ARTIST.NAME
  FROM ARTIST, APPEARS, ALBUM
 WHERE ARTIST.ID=APPEARS.ARTIST_ID
   AND APPEARS.ALBUM_ID=ALBUM.ID
   AND ALBUM.NAME="Mooshoo Tribute"
 ORDER BY ARTIST.ID
```

*Step #1: Decompose into single-value queries*

*Step #2: Substitute the values from*
*        Query#1 → Query #3 → Query #4*

*Query #1*

```
SELECT ALBUM.ID AS ALBUM_ID INTO TEMP1
  FROM ALBUM
 WHERE ALBUM.NAME="Mooshoo Tribute"
```

*Query #3*

```
SELECT APPEARS.ARTIST_ID INTO TEMP2
  FROM APPEARS, TEMP1
 WHERE APPEARS.ALBUM_ID=TEMP1.ALBUM_ID
 ORDER BY APPEARS.ARTIST_ID
```

*Query #4*

```
SELECT ARTIST.NAME
  FROM ARTIST, TEMP2
 WHERE ARTIST.ARTIST_ID=TEMP2.ARTIST_ID
```

# INGRES OPTIMIZER

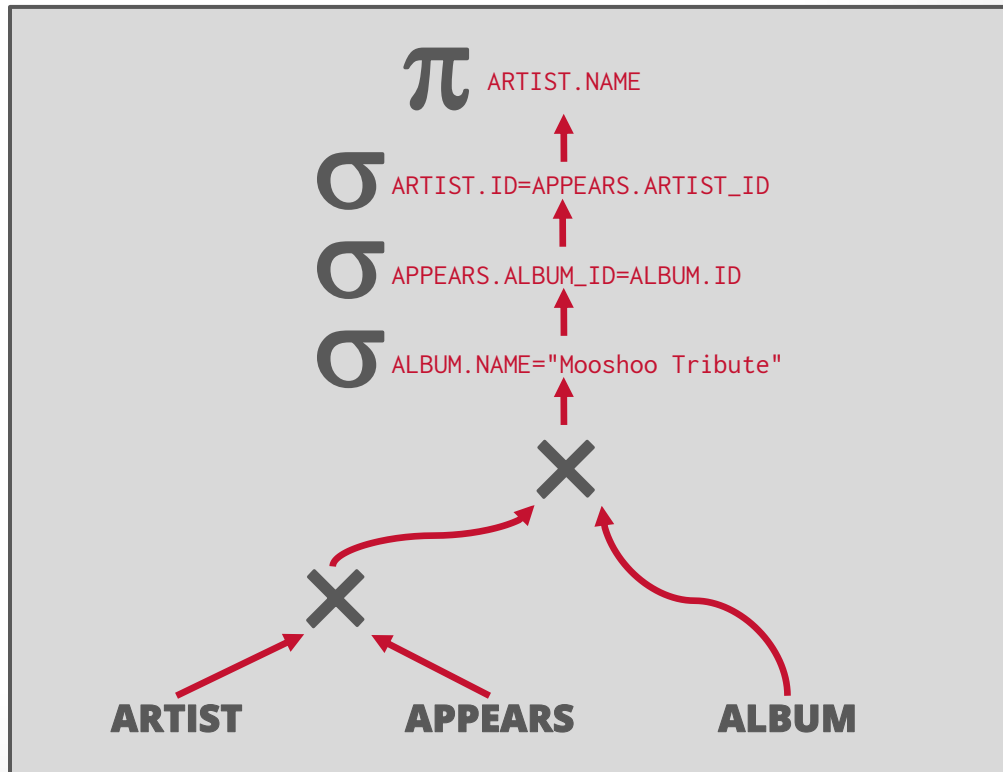*Retrieve the names of people that appear on the DJ Mooshoo Tribute mixtape ordered by their artist id.*
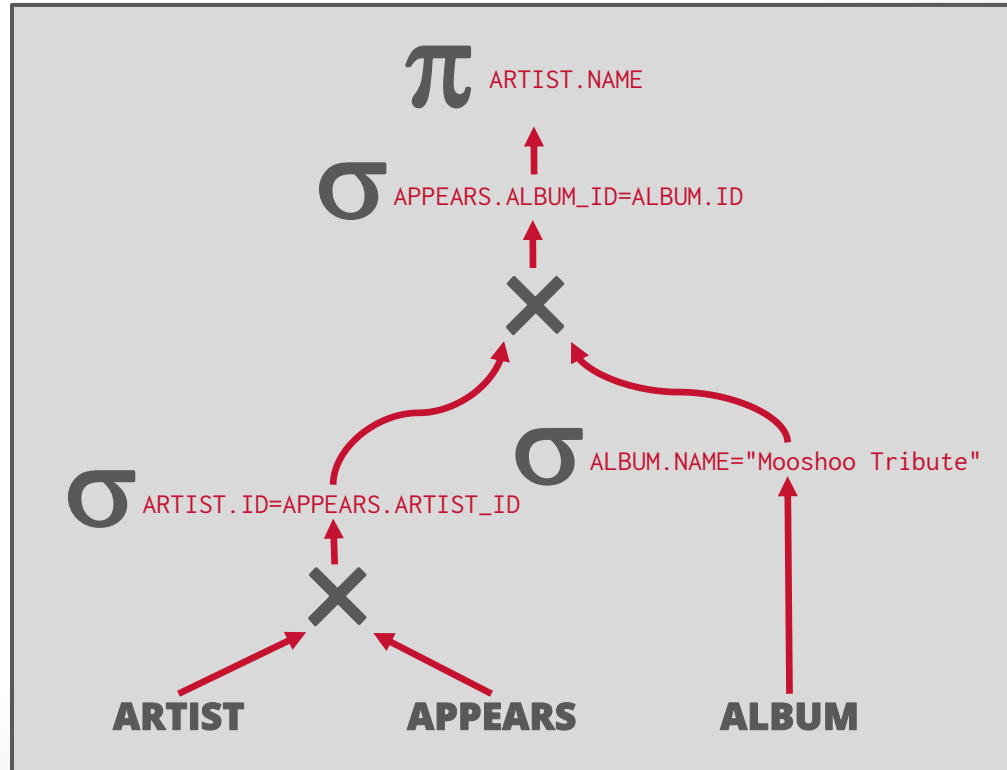
```
SELECT ARTIST.NAME
  FROM ARTIST, APPEARS, ALBUM
 WHERE ARTIST.ID=APPEARS.ARTIST_ID
   AND APPEARS.ALBUM_ID=ALBUM.ID
   AND ALBUM.NAME="Mooshoo Tribute"
ORDER BY ARTIST.ID
```

*Step #1: Decompose into single-value queries*

*Step #2: Substitute the values from Query#1 → Query #3 → Query #4*

| ALBUM_ID |
|----------|
| 9999 |

```
SELECT APPEARS.ARTIST_ID
  FROM APPEARS
 WHERE APPEARS.ALBUM_ID=9999
 ORDER BY APPEARS.ARTIST_ID
```

*Query #4*

```
SELECT ARTIST.NAME
  FROM ARTIST, TEMP2
 WHERE ARTIST.ARTIST_ID=TEMP2.ARTIST_ID
```

# INGRES OPTIMIZER

*Retrieve the names of people that appear on the DJ Mooshoo Tribute mixtape ordered by their artist id.*

```
SELECT ARTIST.NAME
  FROM ARTIST, APPEARS, ALBUM
 WHERE ARTIST.ID=APPEARS.ARTIST_ID
   AND APPEARS.ALBUM_ID=ALBUM.ID
   AND ALBUM.NAME="Mooshoo Tribute"
 ORDER BY ARTIST.ID
```
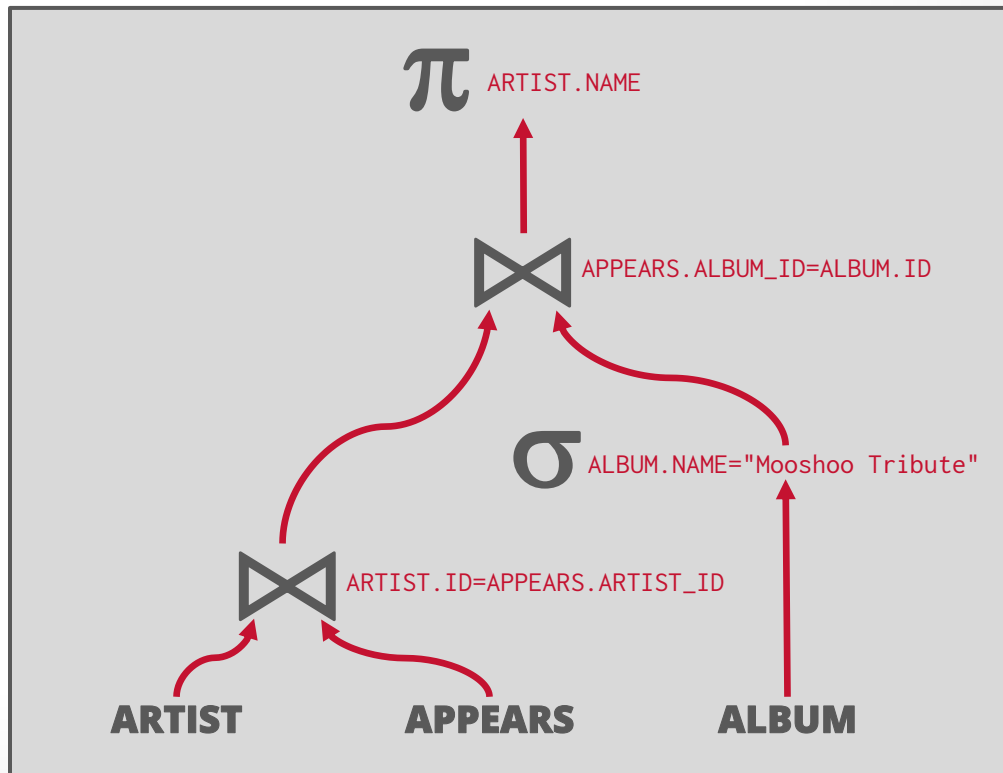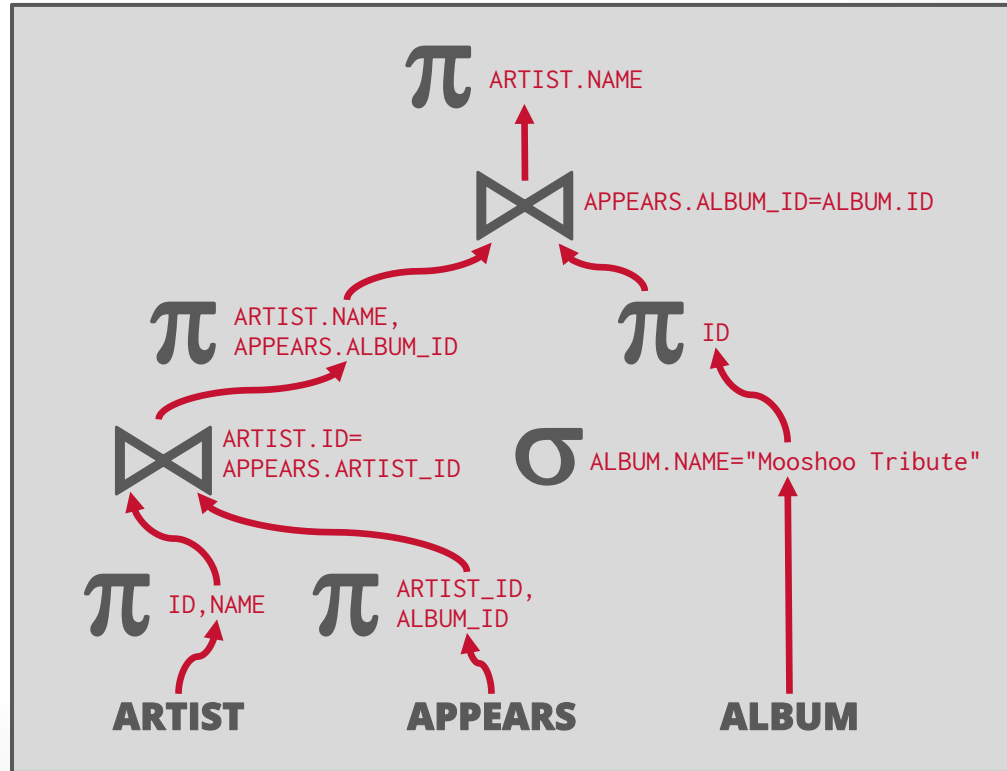
*Step #1: Decompose into single-value queries*

*Step #2: Substitute the values from Query#1 → Query #3 → Query #4*

| ALBUM_ID |
|----------|
| 9999 |

| ARTIST_ID |
|-----------|
| 123 |
| 456 |

```
SELECT ARTIST.NAME
  FROM ARTIST
 WHERE ARTIST.ARTIST_ID=123
```

```
SELECT ARTIST.NAME
  FROM ARTIST
 WHERE ARTIST.ARTIST_ID=456
```

# INGRES OPTIMIZER

*Retrieve the names of people that appear on the DJ Mooshoo Tribute mixtape ordered by their artist id.*

```
SELECT ARTIST.NAME
  FROM ARTIST, APPEARS, ALBUM
 WHERE ARTIST.ID=APPEARS.ARTIST_ID
   AND APPEARS.ALBUM_ID=ALBUM.ID
   AND ALBUM.NAME="Mooshoo Tribute"
ORDER BY ARTIST.ID
```

*Step #1: Decompose into single-value queries*

*Step #2: Substitute the values from*
*Query#1 → Query #3 → Query #4*

| ALBUM_ID |
|----------|
| 9999 |

| ARTIST_ID |
|-----------|
| 123 |
| 456 |

| NAME |
|------|
| O.D.B. |

| NAME |
|------|
| DJ Premier |

# HEURISTIC-BASED OPTIMIZATION

**Advantages:**

→ Easy to implement and debug.

→ Works reasonably well and is fast for simple queries.

**Disadvantages:**

→ Relies on magic constants that predict the efficacy of a planning decision.

→ Nearly impossible to generate good plans when operators have complex inter-dependencies.

# HEURISTIC-BASED OPTIMIZATION

**Advantages:**

→ Easy to implement and debug.

→ Works reasonably well and is fast

**Disadvantages:**

→ Relies on magic constants that pre
planning decision.

→ Nearly impossible to generate goo
have complex inter-dependencies

Stonebraker gave the story of the query optimizer as an example. Relational queries were often highly complex. Let's say you wanted your database to give you the name, salary, and job title of everyone in your Chicago office who did the same kind of work as an employee named Alien. (This example happens to come from Oracle's 1981 user guide.) This would require the database to find information in the employee table and the department table, then sort the data. How quickly the database management system did this depended on how cleverly the system was constructed. "If you do it smart, you get the answer a lot quicker than if you do it stupid, Stonebraker said.

He continued. "Oracle had a really stupid optimizer. They did the query in the order that you happened to type in the clauses. Basically, they blindly did it from left to right. The Ingres program looked at everything there and tried to figure out the best way to do it." But Ellison found a way to neutralize this advantage, Stonebraker said. "Oracle was really shrewd. They said they had a syntactic optimizer, whereas the other guys had a semantic optimizer. The truth was, they had no optimizer and the other guys had an optimizer. It was very, very, very creative marketing. . . . They were very good at confusing the market."

"What he's using is semantics himself," Ellison said. Just because Oracle did things differently, "Stonebraker decided we didn't have an optimizer. [He seemed to think] the only kind of optimizer was his optimizer, and our approach to optimization wasn't really optimization at all. That's an interesting notion, but I'm not sure I buy that."

# HISTORY OF QUERY OPTIMIZERS

**Choice #1: Heuristics**
→ INGRES (1970s), Oracle (until mid 1990s)

**Choice #2: Heuristics + Cost-based Join Search**
→ System R (1970s), early IBM DB2

**Choice #3: Stratified Search**
→ IBM STARBURST (late 1980s), now IBM DB2 + Oracle

**Choice #4: Unified Search**
→ Volcano/Cascades in 1990s, now MSSQL + Greenplum

**Choice #5: Randomized Search**
→ Academics in the 1980s, current Postgres

# HEURISTICS + COST-BASED SEARCH

First evaluate static rules to perform initial logical→logical optimizations.

Then enumerate plans using logical→physical transformations to find best plan according to a cost model.

*Selinger*

**Examples**: System R, early IBM DB2, most open-source DBMSs today.

ACCESS PATH SELECTION IN A RELATIONAL DATABASE MANAGEMENT SYSTEM
*SIGMOD 1979*

# PHYSICAL QUERY OPTIMIZATION

Transform a query plan's logical operators into physical operators.
→ Add more execution information
→ Select indexes / access paths
→ Choose operator implementations
→ Choose when to materialize (i.e., temp tables).

This stage must support cost model estimates.

# SYSTEM R OPTIMIZER

Break query up into blocks and generate the logical operators for each block.

For each logical operator, generate a set of physical operators that implement it.
→ All combinations of join algorithms and access paths

If a block accesses multiple relations, iteratively construct a join tree that minimizes the estimated amount of work to execute the plan.

# SYSTEM R – SINGLE RELATION QUERIES

Access path selection for a single relation query block is (relatively) easy because they are **sargable**.

*Search Argument Able*

Pick the best access method (sequential scan vs. index) using a simple cost model.

```
SELECT id
  FROM xxx
 WHERE val >= 123
   AND val <= 456;
```

$\pi_{id}$

$\sigma_{val=123}$

**xxx**

# SYSTEM R – SINGLE RELATION QUERIES

Access path selection for a single relation query block is (relatively) easy because they are **sargable**.

*Search Argument Able*

Pick the best access method (sequential scan vs. index) using a simple cost model.

```
SELECT id
  FROM xxx
 WHERE val >= 123
   AND val <= 456;
```

$\pi_{id}$

$\sigma_{val=123}$

**XXX**

```
CREATE TABLE xxx (
  id INT PRIMARY KEY,
  val INT,
  ⋮
);
CREATE INDEX ON xxx (val);
```

# SYSTEM R – COST MODEL

The cost of an access method is the summation of the expected number of I/Os ("page fetches") and weighted computational cost ("RSI calls").
→ Weight determines relative cost of I/O versus CPU.

The DBMS estimates these values based on the **selectivity factor** of predicates derived from statistics for each relation and its indexes.

The OPTIMIZER examines both the predicates in the query and the access paths available on the relations referenced by the query, and formulates a cost prediction for each access plan, using the following cost formula:

COST = PAGE FETCHES + W * (RSI CALLS).

This cost is a weighted measure of I/O (pages fetched) and CPU utilization (instructions executed). W is an adjustable weighting factor between I/O and CPU. RSI CALLS is the predicted number of tuples returned from the RSS. Since most of System R's CPU time is spent in the RSS, the number of RSI calls is a good approximation for CPU utilization. Thus the choice of a minimum cost path to process a query attempts to minimize total resources required.

| TABLE 2 SITUATION | COST FORMULAS COST (in pages) |
|---|---|
| Unique index matching an equal predicate | 1 + 1 + W |
| Clustered index I matching one or more boolean factors | F(preds) * (NINDX(I) + TCARD) + W * RSICARD |
| Non-clustered index I matching one or more boolean factors | F(preds) * (NINDX(I) + NCARD) + W * RSICARD |
| | or F(preds) * (NINDX(I) + TCARD) + W * RSICARD if this number fits in the System R buffer |
| Clustered index I not matching any boolean factors | (NINDX(I) + TCARD) + W * RSICARD |
| Non-clustered index I not matching any boolean factors | (NINDX(I) + NCARD) + W * RSICARD |
| | or (NINDX(I) + TCARD) + W * RSICARD if this number fits in the System R buffer |
| Segment scan | TCARD/P + W * RSICARD |

# SYSTEM R – COST MODEL

The cost of an access method is summation of the expected num I/Os ("page fetches") and weigh computational cost ("RSI calls").
→ Weight determines relative cost of versus CPU.

The DBMS estimates these valu based on the **selectivity factor** predicates derived from statistic each relation and its indexes.

## 17.6.2. Planner Cost Constants

**Note:** Unfortunately, there is no well-defined method for determining ideal values for the family of "cost" variables that appear below. You are encouraged to experiment and share your findings.

`random_page_cost` (floating point)

Sets the planner's estimate of the cost of a nonsequentially fetched disk page. This is measured as a multiple of the cost of a sequential page fetch. A higher value makes it more likely a sequential scan will be used, a lower value makes it more likely an index scan will be used. The default is four.

`cpu_tuple_cost` (floating point)

Sets the planner's estimate of the cost of processing each row during a query. This is measured as a fraction of the cost of a sequential page fetch. The default is 0.01.

`cpu_index_tuple_cost` (floating point)

Sets the planner's estimate of the cost of processing each index row during an index scan. This is measured as a fraction of the cost of a sequential page fetch. The default is 0.001.

`cpu_operator_cost` (floating point)

Sets the planner's estimate of the cost of processing each operator in a WHERE clause. This is measured as a fraction of the cost of a sequential page fetch. The default is 0.0025.

| Non-clustered index I not matching any boolean factors | (NINDX(I) + NCARD) + W * RSICARD |
| | or (NINDX(I) + TCARD) + W * RSICARD if this number fits in the System R buffer |
| Segment scan | TCARD/P + W * RSICARD |

# SYSTEM R – SELECTIVITY FACTOR

A selectivity factor of a predicate is the expected faction of tuples that will satisfy that predicate.

The optimizer uses formulas to approximate each predicate's selectivity factor.
→ Make several assumptions about distribution of values in columns to simplify the problem.

# SYSTEM R – SELECTIVITY FACTOR

A selectivity factor of a predicate is the expected faction of tuples that will satisfy that predicate.

The optimizer u...
approximate eac...
selectivity factor.
→ Make several assumptions about distribution of values in columns to simplify the problem.



$$F = 1 / \text{ICARD(column index)} \text{ if there is an index on column}$$
This assumes an even distribution of tuples among the index key values.
$$F = 1/10 \quad \text{otherwise}$$

# SYSTEM R – INTERESTING ORDERS

For each query block, the DBMS extracts the required ("interesting") ordering of its output.
→ Examples: **ORDER BY**, **GROUP BY**

It then compares the best access method that orders the data versus the best unordered access method + sort operator.

If there is no required ordering, then the DBMS selects the access method with the lowest cost.

# SYSTEM R — MULTIPLE RELATIONS

If a query block accesses multiple relations, then the DBMS must determine the best ordering to join those relations.
→ Also identify interesting orders based on join predicates.

Leverage domain knowledge to reduce the search complexity by delaying or discarding plan choices.
→ Example: Only consider left-deep trees.

Join costs are estimated based on the number of tuples processed in outer/inner relations.

# SYSTEM R – MULTIPLE RELATIONS

*Step #1: Choose the best access paths to each relation.*

*Step #2: Enumerate all join orderings for 1-relation plans using best access path found in Step #1.*

*Step #3: For each subsequent pass, the algorithm determines the best way to join the result of an n − 1 relation plan as the outer relation to the nth relation.*

Algorithm does <u>not</u> need to remember anything at a previous level explicitly as it's being remembered implicitly by the nature of a bottom-up approach.

# SYSTEM R – MULTIPLE RELATIONS

*Retrieve the names of people that appear on Andy's mixtape ordered by their artist id.*

```
SELECT ARTIST.NAME
  FROM ARTIST, APPEARS, ALBUM
 WHERE ARTIST.ID=APPEARS.ARTIST_ID
   AND APPEARS.ALBUM_ID=ALBUM.ID
   AND ALBUM.NAME="Mooshoo Tribute"
 ORDER BY ARTIST.ID
```

**ARTIST**: Sequential Scan

**APPEARS**: Sequential Scan

**ALBUM**: Index Look-up on **NAME**

*Step #1: Choose the best access paths to each table*

# SYSTEM R – MULTIPLE RELATIONS

*Retrieve the names of people that appear on Andy's mixtape ordered by their artist id.*

```
SELECT ARTIST.NAME
  FROM ARTIST, APPEARS, ALBUM
 WHERE ARTIST.ID=APPEARS.ARTIST_ID
   AND APPEARS.ALBUM_ID=ALBUM.ID
   AND ALBUM.NAME="Mooshoo Tribute"
ORDER BY ARTIST.ID
```

*Step #1: Choose the best access paths to each table*

*Step #2: Enumerate all possible join orderings for tables*

**ARTIST**: Sequential Scan

**APPEARS**: Sequential Scan

**ALBUM**: Index Look-up on **NAME**

ARTIST  ⋈  APPEARS  ⋈  ALBUM

APPEARS  ⋈  ALBUM   ⋈  ARTIST

ALBUM   ⋈  APPEARS  ⋈  ARTIST

APPEARS  ⋈  ARTIST  ⋈  ALBUM

ARTIST  ✕  ALBUM   ⋈  APPEARS

ALBUM   ✕  ARTIST  ⋈  APPEARS

⋮              ⋮              ⋮

# SYSTEM R – MULTIPLE RELATIONS

*Retrieve the names of people that appear on Andy's mixtape ordered by their artist id.*

```
SELECT ARTIST.NAME
  FROM ARTIST, APPEARS, ALBUM
 WHERE ARTIST.ID=APPEARS.ARTIST_ID
   AND APPEARS.ALBUM_ID=ALBUM.ID
   AND ALBUM.NAME="Mooshoo Tribute"
ORDER BY ARTIST.ID
```

**ARTIST**: Sequential Scan
**APPEARS**: Sequential Scan
**ALBUM**: Index Look-up on **NAME**

*Step #1: Choose the best access paths to each table*

*Step #2: Enumerate all possible join orderings for tables*

*Step #3: Determine the join ordering with the lowest cost*

ARTIST   ⋈   APPEARS   ⋈   ALBUM

APPEARS   ⋈   ALBUM   ⋈   ARTIST

ALBUM   ⋈   APPEARS   ⋈   ARTIST

APPEARS   ⋈   ARTIST   ⋈   ALBUM

ARTIST   ✕   ALBUM   ⋈   APPEARS

ALBUM   ✕   ARTIST   ⋈   APPEARS

⋮        ⋮        ⋮

Logical Op
Physical Op

# SYSTEM R - BOTTOM-UP SEARCH

ARTIST ⋈ APPEARS ⋈ ALBUM



ARTIST⋈APPEARS
ALBUM

ALBUM⋈APPEARS
ARTIST

APPEARS⋈ALBUM
ARTIST

• • •

NL_JOIN(A1,A3)    MERGE_JOIN(A1,A3)    NL_JOIN(A2,A3)    MERGE_JOIN(A2,A3)    NL_JOIN(A3,A2)    MERGE_JOIN(A3,A2)    • • •

ALBUM.ID=APPEARS.ALBUM_ID

ARTIST.ID=APPEARS.ARTIST_ID

APPEARS.ALBUM_ID=ALBUM.ID

ARTIST ALBUM APPEARS

*Logical Op*
*Physical Op*

# SYSTEM R – BOTTOM-UP SEARCH

ARTIST ⋈ APPEARS ⋈ ALBUM

ARTIST⋈APPEARS
ALBUM

ALBUM⋈APPEARS
ARTIST

APPEARS⋈ALBUM
ARTIST

● ● ●

NL_JOIN(A1,A3)

NL_JOIN(A2,A3)

MERGE_JOIN(A3,A2)

● ● ●

ALBUM.ID=APPEARS.ALBUM_ID

ARTIST.ID=APPEARS.ARTIST_ID

APPEARS.ALBUM_ID=ALBUM.ID

ARTIST ALBUM APPEARS

# SYSTEM R – BOTTOM-UP SEARCH

# SYSTEM R — BOTTOM-UP SEARCH

Logical Op

Physical Op

ARTIST ⋈ APPEARS ⋈ ALBUM

NL_JOIN(A1⋈A3,A2)

NL_JOIN(A2⋈A3,A1)

NL_JOIN(A3⋈A2,A1)

· · ·

APPEARS.ALBUM_ID=ALBUM.ID

ARTIST⋈APPEARS
ALBUM

APPEARS.ARTIST_ID=ARTIST.ID

ALBUM⋈APPEARS
ARTIST

APPEARS.ARTIST_ID=ARTIST.ID

APPEARS⋈ALBUM
ARTIST

· · ·

NL_JOIN(A1,A3)

ALBUM.ID=APPEARS.ALBUM_ID

NL_JOIN(A2,A3)

MERGE_JOIN(A3,A2)

· · ·

ARTIST.ID=APPEARS.ARTIST_ID

APPEARS.ALBUM_ID=ALBUM.ID

ARTIST ALBUM APPEARS

Logical Op

*Physical Op*

# SYSTEM R – BOTTOM-UP SEARCH

ARTIST ⋈ APPEARS ⋈ ALBUM

NL_JOIN(A2⋈A3,A1)

APPEARS.ARTIST_ID=ARTIST.ID

ALBUM⋈APPEARS
ARTIST

NL_JOIN(A2,A3)

ALBUM.ID=APPEARS.ALBUM_ID

ARTIST ALBUM APPEARS

# SYSTEM R – BOTTOM-UP SEARCH

ARTIST ⋈ APPEARS ⋈ ALBUM

NL_JOIN(A2⋈A3,A1)

*The query has* **ORDER BY** *on* **ARTIST.ID** *but the logical plans do not contain sorting properties.*

APPEARS.ARTIST_ID=ARTIST.ID

ALBUM⋈APPEARS
ARTIST

NL_JOIN(A2,A3)

*Keep track of best plans with and without data in proper physical form, and then check whether tacking on a sort operator at the end is better.*

ALBUM.ID=APPEARS.ALBUM_ID

ARTIST ALBUM APPEARS

# PLAN ENUMERATION

## Approach #1: Generative / Bottom-Up
→ Start with nothing and then iteratively assemble and add building blocks to generate a query plan.
→ **Examples:** System R, Starburst

## Approach #2: Transformation / Top-Down
→ Start with the outcome that the query wants and then transform it to equivalent alternative sub-plans to find the optimal plan that gets to that goal.
→ **Examples**: Volcano, Cascades

# SYSTEM R — NESTED QUERIES

The DBMS treats nested queries as separate queries.

```
SELECT name FROM employee
 WHERE salary >(SELECT AVG(salary)
                FROM employee);
```

The optimizer executes an inner query before it begins planning an outer query so that it can substitute values into it or materialize its results to a temporary table.

We will have an entire lecture on rewriting nested queries into joins…

# SYSTEM R – NESTED QUERIES

The DBMS treats nested queries as separate queries.

The optimizer executes an inner query before it begins planning an outer query so that it can substitute values into it or materialize its results to a temporary table.

We will have an entire lecture on rewriting nested queries into joins…

```
SELECT name FROM employee
 WHERE salary >(SELECT AVG(salary)
                FROM employee);
```

```
SELECT AVG(salary) FROM employee;
```

# SYSTEM R — NESTED QUERIES

The DBMS treats nested queries as separate queries.

The optimizer executes an inner query before it begins planning an outer query so that it can substitute values into it or materialize its results to a temporary table.

We will have an entire lecture on rewriting nested queries into joins…

```
SELECT name FROM employee
 WHERE salary >(SELECT AVG(salary)
                FROM employee);
```

```
SELECT AVG(salary) FROM
```

| AVG(salary) |
| --- |
| 100000 |

# SYSTEM R – NESTED QUERIES

The DBMS treats nested queries as separate queries.

The optimizer executes an inner query before it begins planning an outer query so that it can substitute values into it or materialize its results to a temporary table.

We will have an entire lecture on rewriting nested queries into joins…

```
SELECT name FROM employee
 WHERE salary >(SELECT AVG(salary)
               FROM employee);
```

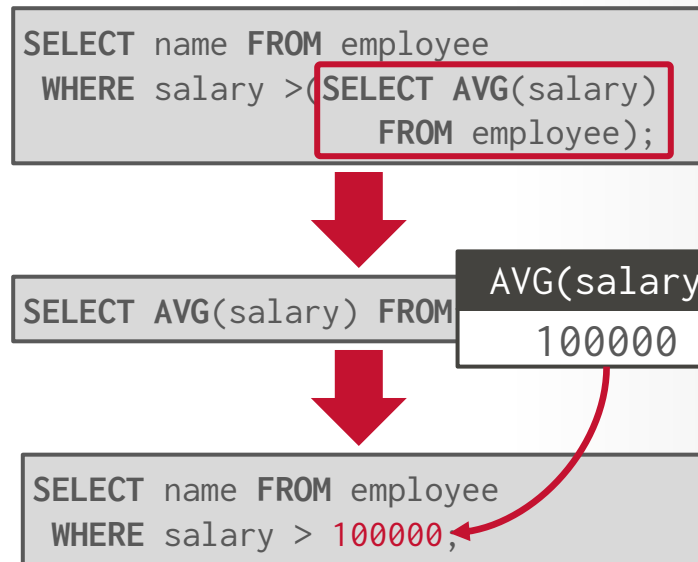| AVG(salary) |
|-------------|
| 100000 |

```
SELECT AVG(salary) FROM
```

```
SELECT name FROM employee
 WHERE salary > 100000;
```

# HEURISTICS + COST-BASED SEARCH

**Advantages:**

→ Usually finds a reasonable plan without having to perform an exhaustive search.

**Disadvantages:**

→ All the same problems as the heuristic-only approach.
→ Left-deep join trees are not always optimal.
→ Must take in consideration the physical properties of data in the cost model (e.g., sort order).

# PARTING THOUGHTS

Although the System R paper is over 40 years old, it still provides a reasonable foundation for building a modern query optimizer.
→ For two relation queries, it will find the optimal join ordering quickly.

But many of its simplifying assumptions in its cost estimates and selectivity factor cause problems in the real-world.

IBM Starburst Optimizer

# HISTORY OF QUERY OPTIMIZERS

**Choice #1: Heuristics**
→ INGRES (1970s), Oracle (until mid 1990s)

**Choice #2: Heuristics + Cost-based Join Search**
→ System R (1970s), early IBM DB2

**Choice #3: Stratified Search**
→ IBM STARBURST (late 1980s), now IBM DB2 + Oracle

**Choice #4: Unified Search**
→ Volcano/Cascades in 1990s, now MSSQL + Greenplum

**Choice #5: Randomized Search**
→ Academics in the 1980s, current Postgres