# Beyond EXPLAIN

## Query Optimization
## From Theory To Code

Yuto Hayamizu

Ryoji Kawamichi

# Historically …

## Before Relational …

- Querying was **physical**

- Need to understand physical organization
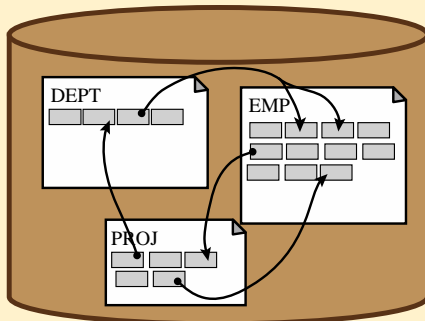
- Navigate query execution by yourself

*"Which file is this table stored in?"*
*"How are records linked?"*
*"Which access path is fast for this table?"*
*"What is the best order of joining tables"*
*…*

# Historically …

## Before Relational …

- Querying was ***physical***

- Need to understand physical organization
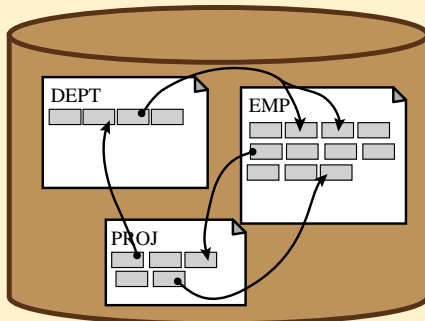
- Navigate query execution by yourself

*"Which file is this table stored in?"*
*"How are records linked?"*
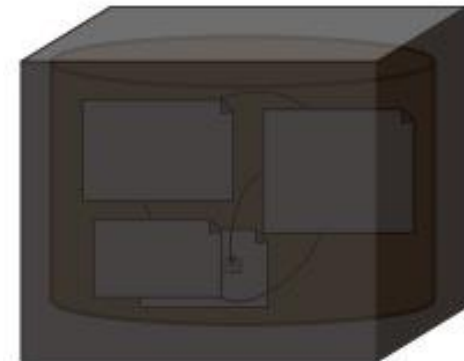*"Which access path is fast for this table?"*
*"What is the best order of joining tables"*

*…*



## After Relational …
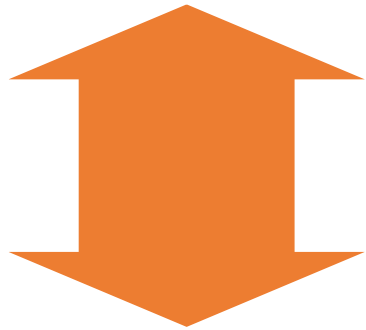
- Querying is logical

- Physical organization is **black-boxed**

- Just declare what you want

# Fill the Gap: *Physical* and Logical

Relational view of data

$$\text{SELECT * FROM } DEPT\, D,\ EMP\, E$$
$$\text{WHERE } E.D\_ID = D.ID \text{ AND } ...$$

## Query Optimizer

- Storage I/O strategy
- Access path selection
- Join method selection
- Aggregation, sorting
- Resource allocation
- ...

# If optimizer perfectly fills the gap…

# We don't need EXPLAIN

# Reality Is Tough

- Optimizer is NOT PERFECT
  - Generated plans are not always optimal, sometimes far from optimal

- We have to take care of physical behavior

- That's why EXPLAIN is so much explained

# Go Beyond EXPLAIN

- Deeper understanding of optimization, better control of your databases


- Theoretical fundamentals of query optimization
  - From basic framework to cutting-edge technologies
- PostgreSQL Optimizer implementation
  - Focusing on basic scan and join methods
  - Behavior observation with TPC-H benchmark

# Outline

- Introduction
- Theory: Query Optimization Framework
- Code: PostgreSQL Optimizer
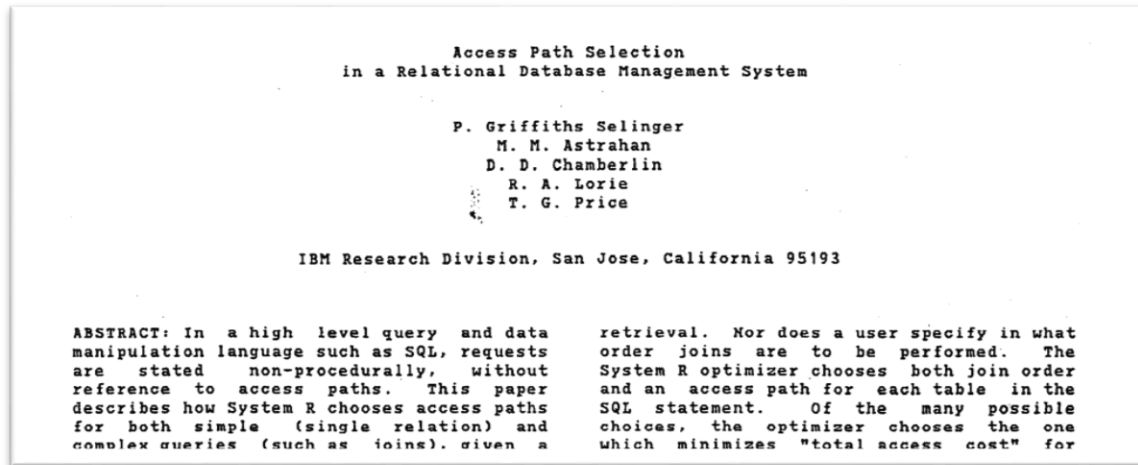- Theory: Cutting-Edge Technologies Overview
- Summary

# Query Optimization Framework

- **Cost-based optimization**
  - Plan selection with estimated execution cost
  - Most of modern optimizers, including PostgreSQL, are cost-based


- Rule-based optimization
  - Plan selection with heuristically ranked rules
  - Easy to produce the same result
  - Hard to evaluate wide variety of plans
  - Ex) Oracle (~10g), Hive (~0.13)

# Main Challenges in Cost-based Optimization

- <u>Cost modeling</u> is HARD
  - Overhead of CPU, I/O, memory access, network, …

- <u>Cardinality estimation</u> is HARD
  - Output size of scans, joins, aggregations, …

- <u>Join ordering search</u> is HARD
  - Combinatorial explosion of join ordering and access path
  - Exhaustive search is NP-hard

# System-R optimizer (1979)

Access Path Selection
in a Relational Database Management System

P. Griffiths Selinger
M. M. Astrahan
D. D. Chamberlin
R. A. Lorie
T. G. Price

IBM Research Division, San Jose, California 95193

ABSTRACT: In a high level query and data manipulation language such as SQL, requests are stated non-procedurally, without reference to access paths. This paper describes how System R chooses access paths for both simple (single relation) and complex queries (such as joins), given a retrieval. Nor does a user specify in what order joins are to be performed. The System R optimizer chooses both join order and an access path for each table in the SQL statement. Of the many possible choices, the optimizer chooses the one which minimizes "total access cost" for

- "The standard"
  - Cost estimation with I/O and CPU
  - Cardinality estimation with table statistics
  - Bottom-up plan search
- Many of modern optimizers are "System-R style"
  - PostgreSQL, MySQL, DB2, Oracle, …

# Cost/Cardinality Estimation

COST = [#page fetched] + W * [#storage API calls]
**I/O cost**    **CPU cost**

weight parameter

- **[#page fetched],[#storage API calls]** are estimated with cost formula and following statistics

  - NCARD(T) … the cardinality of relation T
  - TCARD(T) … the number of pages in relation T
  - ICARD(I) … the number of distinct keys in index I
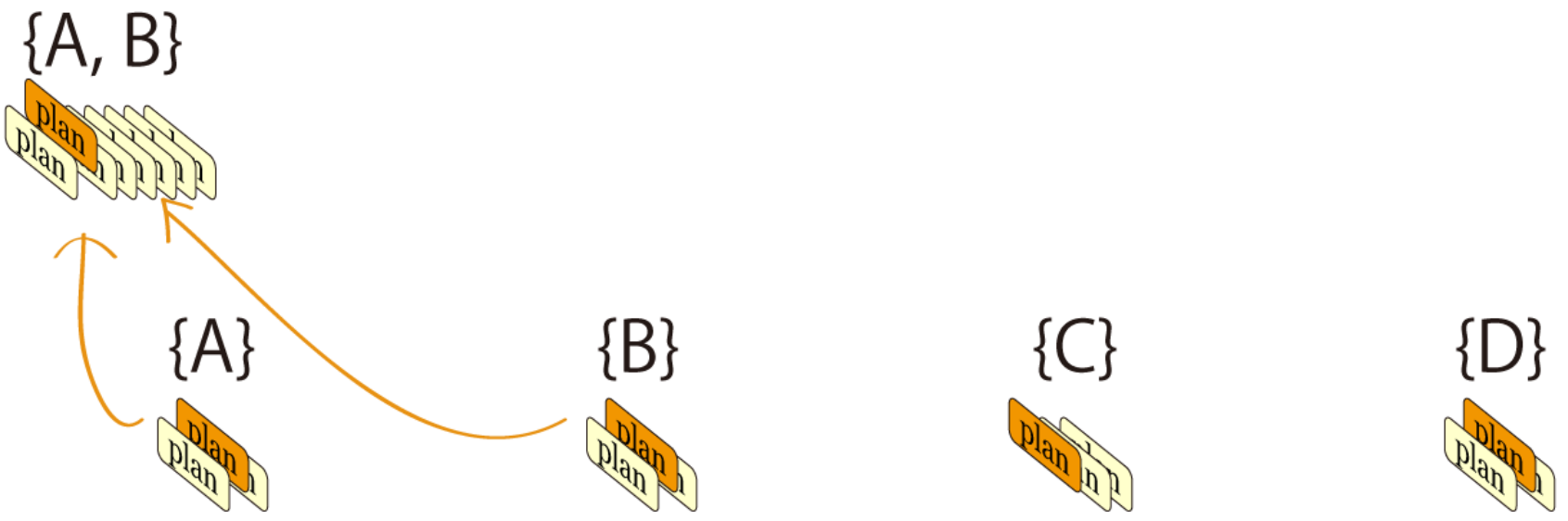  - NINDX(I) … the number of pages in index I

# Bottom-up Plan Search

- Candidate plans for single relation
  - The cheapest access path

- N-relation join ordering search
  - Select the cheapest plans for each relation
  - Then, find optimal join orderings of every 2-relation join
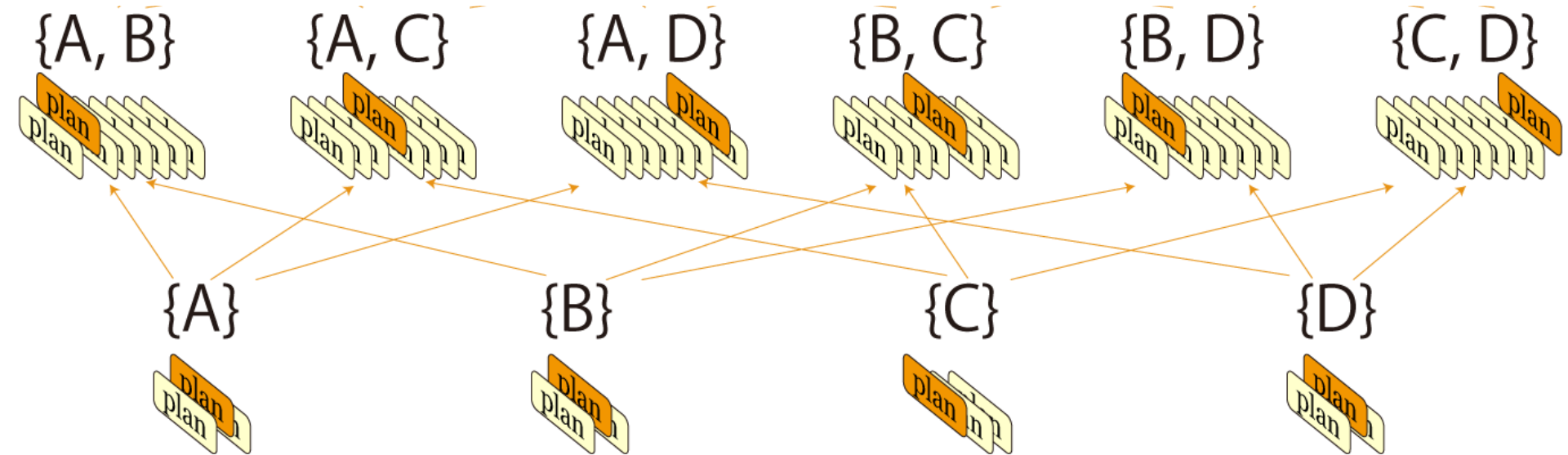  - Then, find optimal join orderings of every 3-relation join
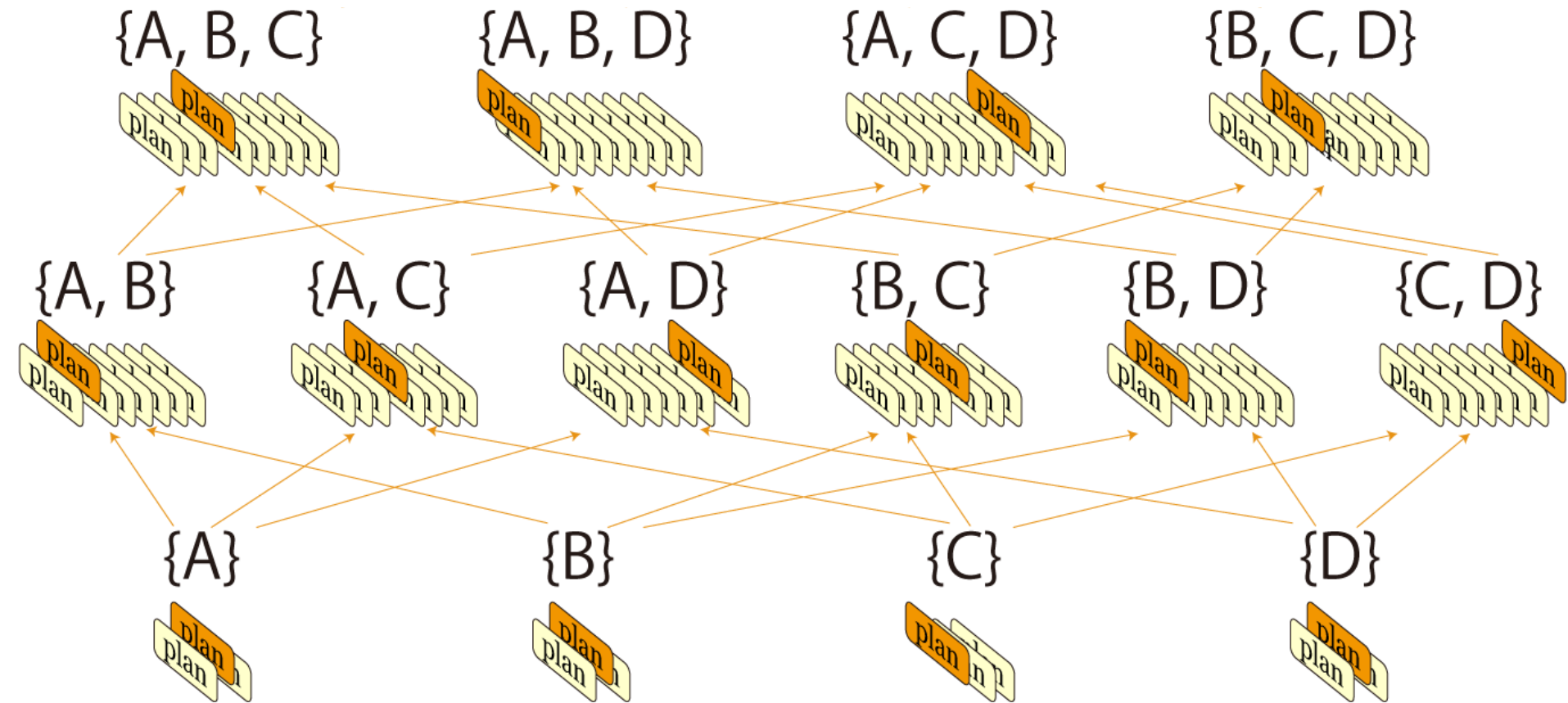    - … until N-relation

Ex) A ⋈ B ⋈ C ⋈ D

{A}          {B}          {C}          {D}

         

Ex) A ⋈ B ⋈ C ⋈ D

{A, B}



{A}       {B}         {C}         {D}

Ex) A ⋈ B ⋈ C ⋈ D

{A, B}    {A, C}    {A, D}    {B, C}    {B, D}    {C, D}

{A}       {B}       {C}       {D}

# Ex) A ⋈ B ⋈ C ⋈ D



{A, B, C}   {A, B, D}   {A, C, D}   {B, C, D}

{A, B}   {A, C}   {A, D}   {B, C}   {B, D}   {C, D}

{A}   {B}   {C}   {D}

# Ex) A ⋈ B ⋈ C ⋈ D

{A, B, C, D}

{A, B, C}    {A, B, D}    {A, C, D}    {B, C, D}

{A, B}    {A, C}    {A, D}    {B, C}    {B, D}    {C, D}

{A}    {B}    {C}    {D}

# Volcano/Cascades (1993)

## The Volcano Optimizer Generator: Extensibility and Efficient Search

**Goetz Graefe**
Portland State University
graefe@cs.pdx.edu

**William J. McKenna**
University of Colorado at Boulder
bill@cs.colorado.edu

**Abstract**

*Emerging database application domains demand not only new functionality but also high performance. To satisfy these two requirements, the Volcano project provides efficient, extensible tools for query and request processing, particularly for object-oriented and scientific database systems. One of these tools is a new optimizer generator. Data model, logical algebra, physical algebra, and optimization rules are translated by the optimizer generator into optimizer source code. Compared with our earlier EX-*

First, this new optimizer generator had to be usable both in the Volcano project with the existing query execution software as well as in other projects as a stand-alone tool. Second, the new system had to be more efficient, both in optimization time and in memory consumption for the search. Third, it had to provide effective, efficient, and extensible support for physical properties such as sort order and compression status. Fourth, it had to permit use of heuristics and data model semantics to guide the search and to prune futile parts of the search space. Finally, it

- Top-down transformational plan search
  - Yet another optimization approach
  - Not well known as "System-R style", but widely used in practice
    Ex) SQL Server, Apache Hive (Apache Calcite), Greenplum Orca
- Extensible optimization framework
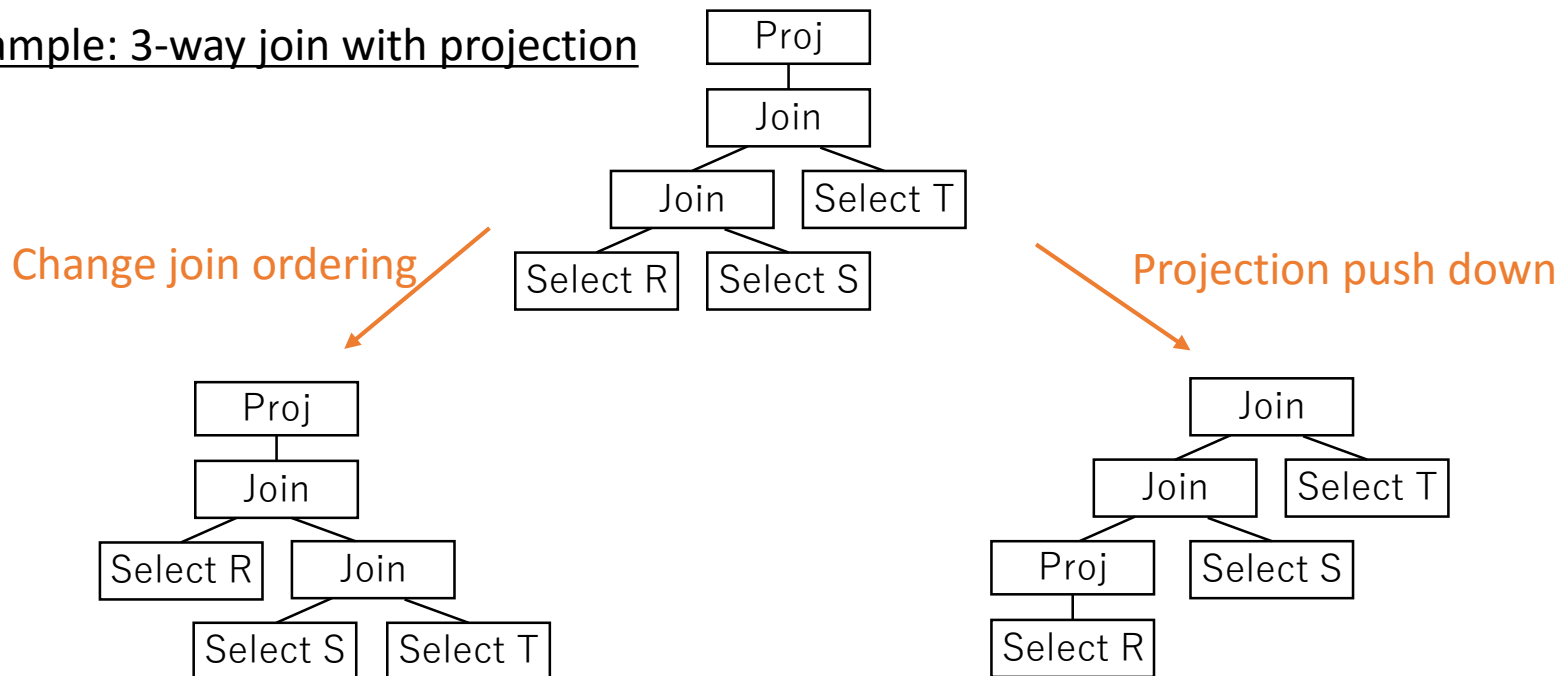
# Extensible Optimization Framework

Query Optimizer <u>Generator</u>

- Generalized expression of query plan not limited to relational data model

- Users (optimizer developers) defines actual implementations:
    - Logical operator ... corresponds to relational algebra
    - Physical algorithm ... corresponds to scan & join methods such as sequential scan, index scan, hash join, nested loop join

# Top-down Transformational Search

- Starts from an initial "logical plan"
- Generate alternative plans with:
  - A) Logical operator transformation
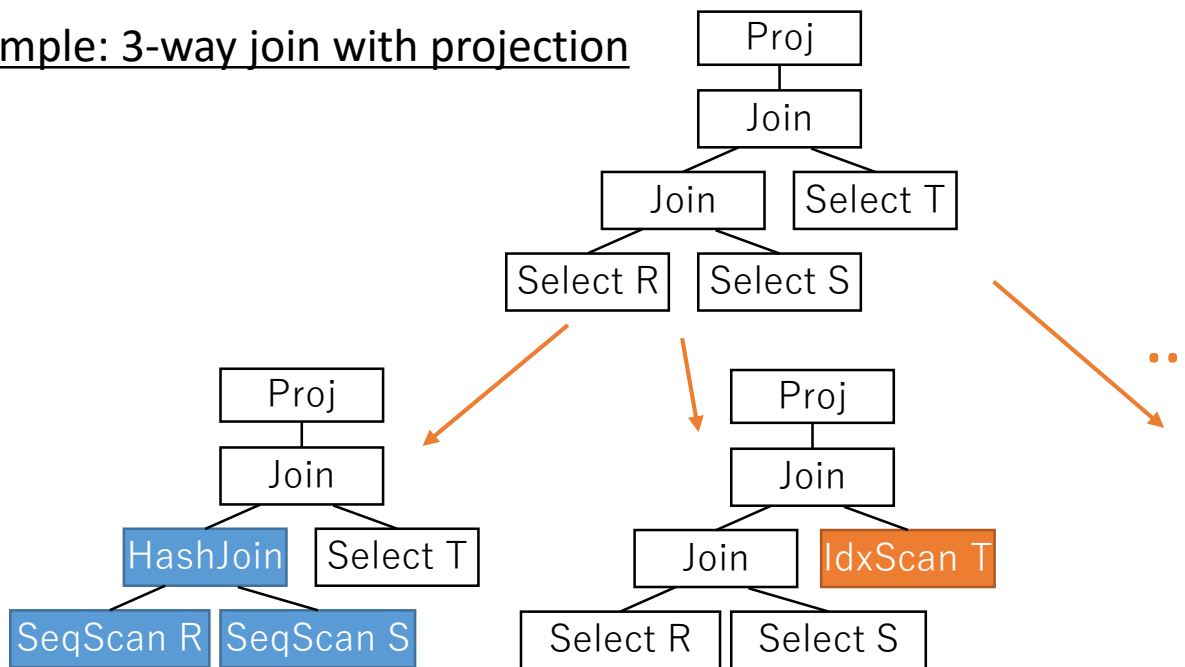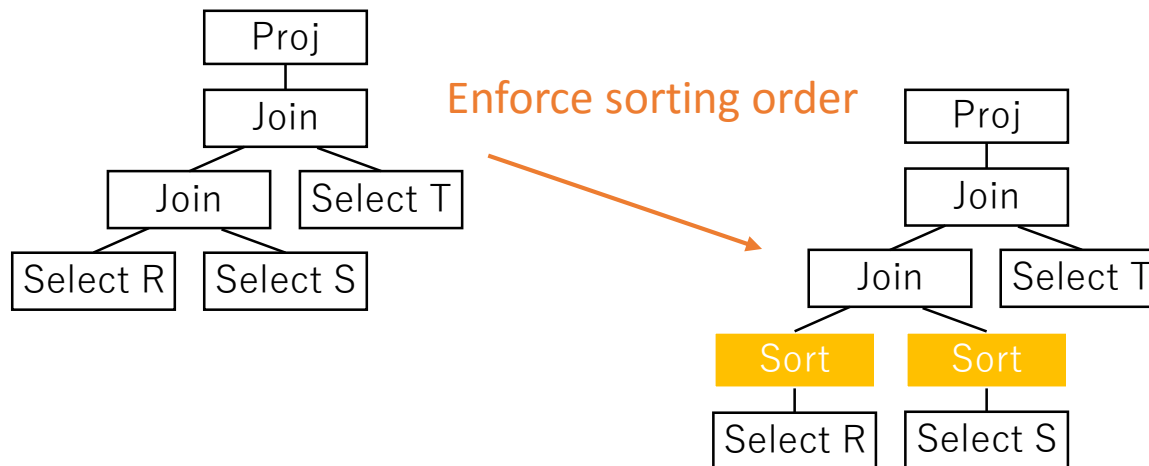  - B) Physical algorithm selection
  - C) Enforcing sorting order

Example: 3-way join with projection



Change join ordering

Projection push down
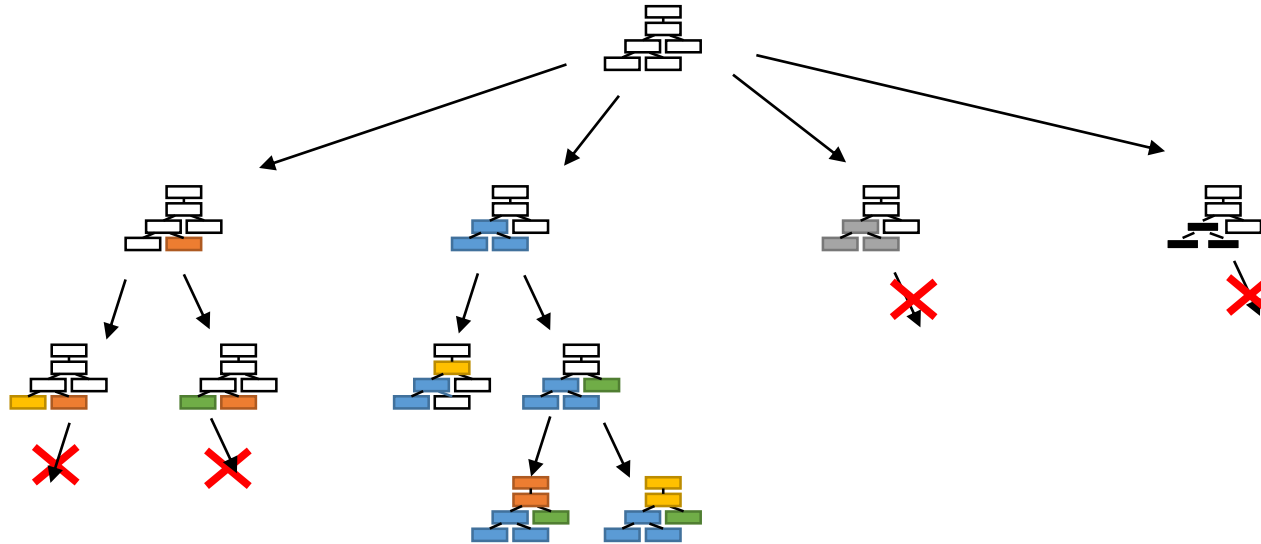
# Top-down Transformational Search

- Starts from an initial "logical plan"
- Generate alternative plans with:
  - A) Logical operator transformation
  - B) Physical algorithm selection
  - C) Enforcing sorting order

Example: 3-way join with projection

# Top-down Transformational Search

- Starts from an initial "logical plan"
- Generate alternative plans with:
  - A) Logical operator transformation
  - B) Physical algorithm selection
  - C) Enforcing sorting order

<u>Example: 3-way join with projection</u>



Enforce sorting order

merge join of R and S is possible now

# Benefits of Top-down approach



- Possible to intentionally limit search space
  - Effective pruning with branch-and-bound
  - Limit search space with search time deadline

# Cost-based Optimization Basics

Two major cost-based optimization style

- System-R
  - Cost modeling with statistics
  - Bottom-up search
- Volcano/Cascades
  - Extensible optimizer generator
    - Cost estimation is user's responsibility
  - Top-down transformational search

# Outline

- Introduction
- Theory: Query Optimization Framework
- Code: PostgreSQL Optimizer
- Theory: Cutting-Edge Technologies Overview
- Summary

# PostgreSQL Optimizer

"System-R style" optimization
- Bottom-up plan search with dynamic programming
- CPU and I/O operation based cost modeling

$$\text{Cost} = c_{\text{seq}} n_{\text{seq}} + c_{\text{rand}} n_{\text{rand}} + c_{\text{tup}} n_{\text{tup}} + \cdots$$

*Seq. I/O*       *Random I/O*      *CPU cost per tuple*

$$= \mathbb{C} \cdot \mathbb{N}$$

$\mathbb{C}$ Cost of single operation
- seq_page_cost
- random_page_cost
- cpu_tuple_cost
- cpu_index_tuple_cost
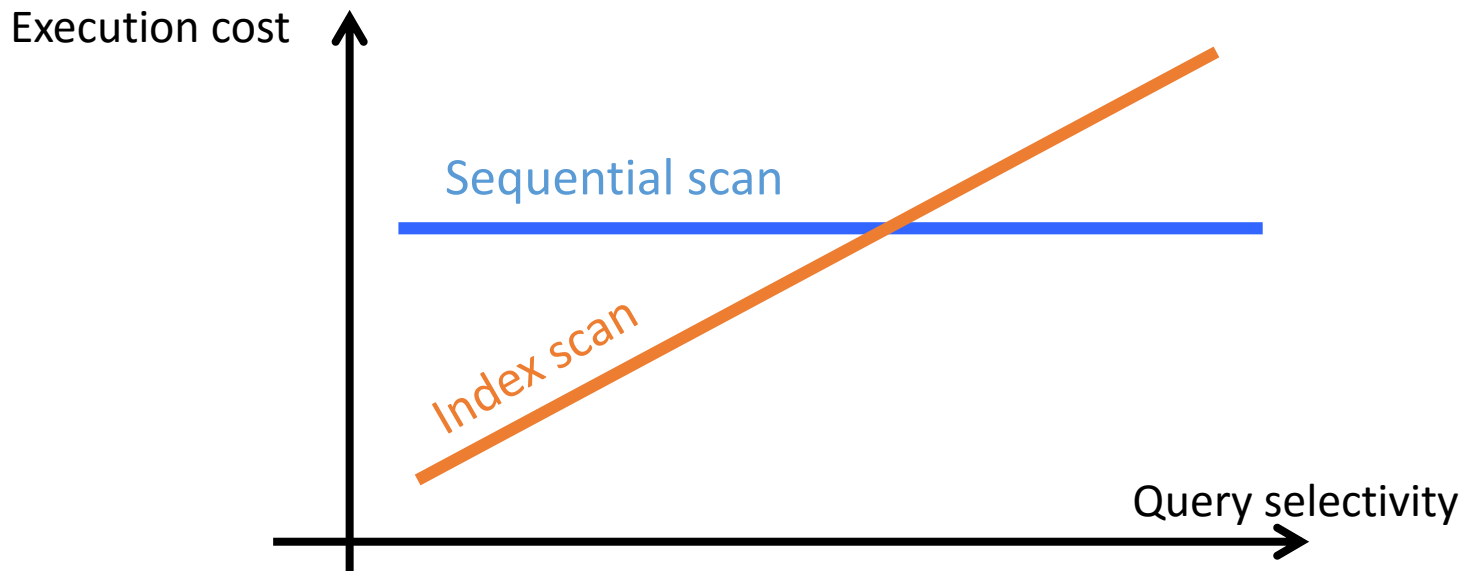- cpu_operator_cost
- (parallel_tuple_cost)

$\mathbb{N}$ Estimated number of each operation
- Cardinality estimation with statistics
- Cost formula for each plan type
  - SeqScan, IndexScan
  - NestLoopJoin, HashJoin, MergeJoin, …

# Detailed Look At Basic Scan Types

- ## Sequential scan
  - Efficient for accessing large potion of tables
- ## Index scan
  - Efficient for accessing a fraction of data

# $\mathbb{N}$ of **SeqScan**

$n_{\text{seq}}$ = (# pages in a table)

$n_{\text{tup}}$ = (# tuples in a table)

$\cdots$ WHERE  A  AND  B  AND $\cdots$

$n_{\text{op}}$ = #qual_operator

= (#tuples) $\times$ (weight factor of A)

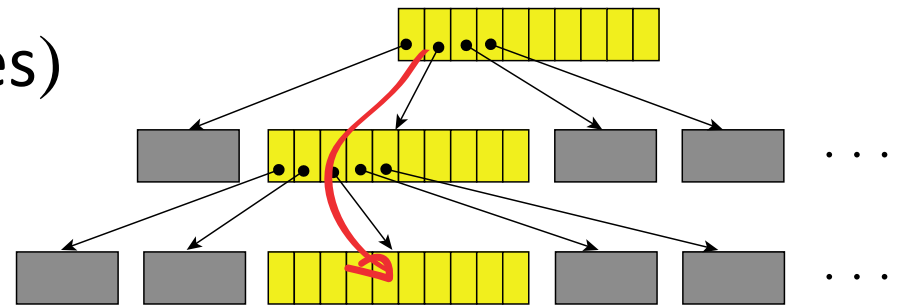$+$ (#tuples) $\times$ (weight factor of B)

$+ \cdots$

# N of **IndexScan**

Consists of:

(A) CPU cost of searching $B^+$-tree

(B) CPU cost of scanning index tuples in leaf pages

(C) I/O cost of leaf pages

(D) I/O cost of heap pages

(E) CPU cost of scanning heap tuples

# 🅝 of **IndexScan**

## (A) B$^+$-tree search

$$n_{\mathrm{op}} \mathrel{+}= \log_2(\text{\#index\_tuples})$$

> I/O cost of internal pages
> Assumed to be always cached in the buffer



## (B) Scanning index tuples in leaf pages

$$n_{\mathrm{itup}} \mathrel{+}= \text{\#qual\_operator}$$
$$\times \ \text{\#leaf\_pages} \times \text{\#ituple\_per\_page} \times \sigma$$

**Selectivity** $\sigma$
Comes from statistics

# $\mathbb{N}$ of **IndexScan**

## (C) I/O cost of index leaf pages

$$n_{\mathrm{rand}} \mathrel{+}= Y(\texttt{effective\_cache\_size}, \#\text{leaf\_pages})$$

**Mackert and Lohman function（Yao function）**

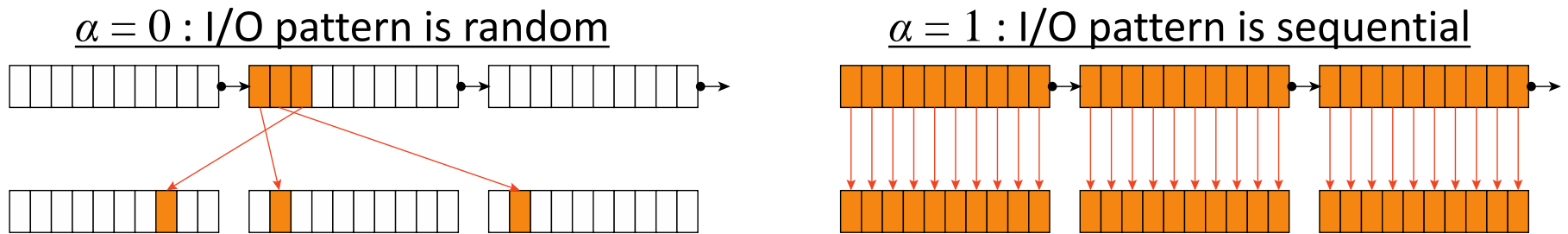I/O count estimation with consideration of buffer caching

$$Y(N,P,\sigma,B) \equiv \begin{cases} \min\left(\dfrac{2PN\sigma}{2P+N\sigma}, P\right) & (P \le B) \\[2ex] \dfrac{2PN\sigma}{2P+N\sigma} & \left(P > B \wedge \sigma \le \dfrac{2PB}{N(2P-B)}\right) \\[2ex] B + \left(N\sigma - \dfrac{2PB}{2P-B}\right)\dfrac{P-B}{P} & \left(P > B \wedge \sigma > \dfrac{2PB}{N(2P-B)}\right) \end{cases}$$



I/O count vs Selectivity σ — no caching effect, Y(N,P,σ,B)

# $\mathbb{N}$ of **IndexScan**

(D) I/O cost of heap pages

Correlation between index and heap ordering: $\alpha$

$\underline{\alpha = 0}$ : I/O pattern is random            $\underline{\alpha = 1}$ : I/O pattern is sequential



$$n_{\mathrm{seq}} \; += \alpha^2 \; \times \; \text{\#match\_pages}$$

$$n_{\mathrm{rand}} += (1\text{-}\alpha^2) \; \times \; \text{\#match\_tuples}$$

(E) CPU cost of scanning heap tuples
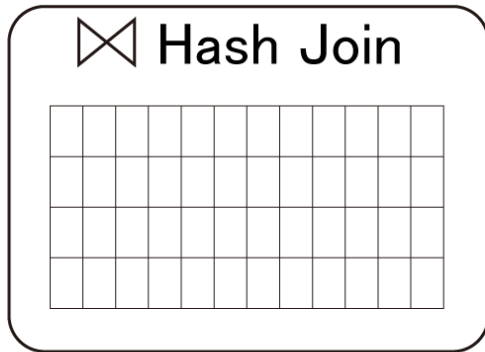
· Estimate the number of scanned tuples from $\boldsymbol{\sigma}$

# Detailed Look At Join Methods

- ## Hash join
  - Efficient for joining large number of records
  - Usually combined with <u>sequential scans</u>

- ## Nested Loop Join
  - Efficient for joining small number of records
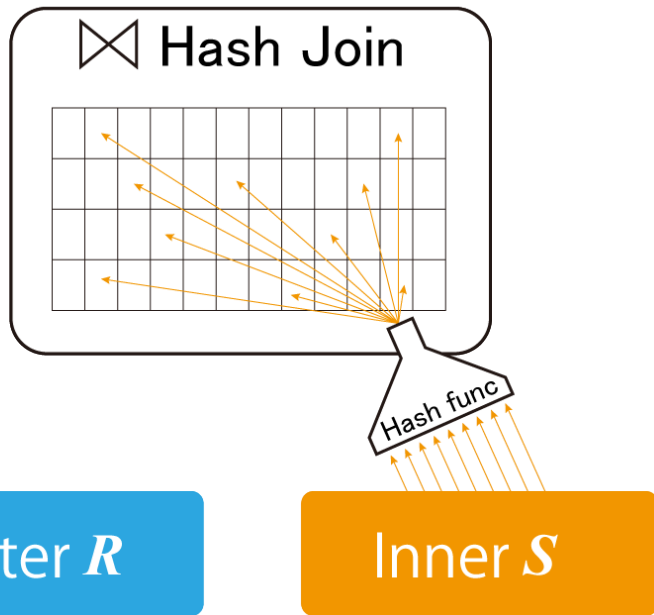  - Usually combined with <u>index scans</u> or <u>small table sequential scans</u>

# N of **HashJoin**

⋈ Hash Join
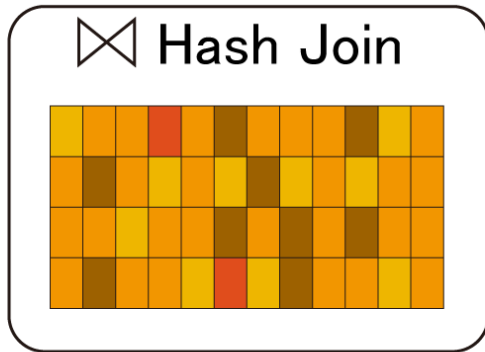
Outer *R*  Inner *S*

# [N] of **HashJoin**

Build phase

- $\text{Cost} += \text{Cost}(\textbf{\textit{inner}})$

  $n_{\text{op}}$ += #qual_op × #inner_tuples

  $n_{\text{tup}}$ += #inner_tuples

  Hashing cost



Hash Join

Hash func

Outer $\textbf{\textit{R}}$

Inner $\textbf{\textit{S}}$

# Ⓝ of **HashJoin**

## Build phase

- $\mathrm{Cost} \mathrel{+}= \mathrm{Cost}(\textbf{\textit{inner}})$

  $n_{\mathrm{op}}$ += #qual_op × #inner_tuples

  $n_{\mathrm{tup}}$ += #inner_tuples

⋈ Hash Join

Outer $\textbf{\textit{R}}$    Inner $\textbf{\textit{S}}$

# $\mathbb{N}$ of **HashJoin**



Hash Join

Outer $R$

Inner $S$

## Build phase

- $\text{Cost} \mathrel{+}= \text{Cost}(\textit{inner}) + \mathbb{C} \cdot \mathbb{N}$

  $n_{\text{op}}$ += #qual_op × #inner_tuples

  $n_{\text{tup}}$ += #inner_tuples

## Probe phase

- $\text{Cost} \mathrel{+}= \text{Cost}(\textit{outer}) + \mathbb{C} \cdot \mathbb{N}$

  $n_{\text{op}}$ += #qual_op × (1 + **#bucket_size × 0.5)**
  × #outer_tuples

  Hashing & table lookup (bucket search) cost

  $n_{\text{tup}}$ += #match_tuples

# N of **HashJoin**

**16 records**

build

#buckets: 2

4 tuples are compared for lookup in average

#buckets : 4

2 tuples are compared for lookup in average

tuple

## Estimated cost of 2-way HashJoin



# of records

# of **NestLoopJoin**

$$R \bowtie S$$

*outer*    *inner*

# $\mathbb{N}$ of **NestLoopJoin**

$$\underset{outer}{\mathbf{R}} \bowtie \underset{inner}{\mathbf{S}}$$



- When #outer_tuples = 1

$$\text{Cost} = \text{Cost}(outer) + \text{Cost}(inner) + \mathbb{C} \cdot \mathbb{N}$$

$n_{\text{tup}}$ += #inner_tuples

$n_{\text{op}}$ += #qual_operator × #inner_tuples

# $\mathbb{N}$ of **NestLoopJoin**

$$R \bowtie S$$

*outer*    *inner*



- When #outer_tuples > 1

$$\text{Cost} = \text{Cost}(outer) + \text{Cost}(inner) + \mathbb{C} \cdot \mathbb{N}$$
$$+ \text{(\#outer\_tuples - 1)} \times \text{Cost}(\textbf{ReScan } inner)$$

Higher buffer hit ratio in ReScan
→ Cost of ReScan is lower than cost of IndexScan

$n_{\text{tup}}$ += #inner_tuples × #outer_tuples

$n_{\text{op}}$ += #qual_operator × #inner_tuples × #outer_tuples

# See How It Works

- TPC-H Benchmark
  - Specification and tools for benchmarking data warehouse workload
    - Open source implementation: DBT-3, pg_tpch
  - Schema, data generation rules and queries
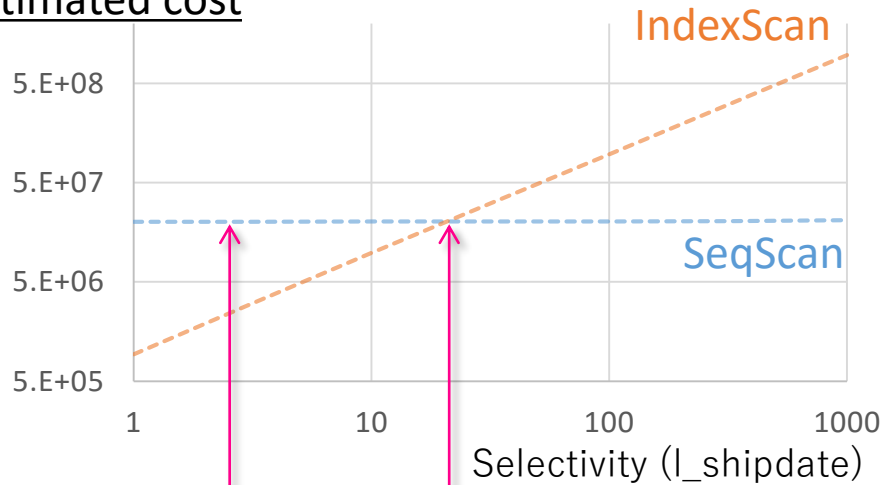
- Experiments with 100GB
  - Scale Factor = 100

# Experimental Setup

- Dell R720xd
    - Xeon (2sockets, 16cores)
    - x24 NL-SAS HDD


- With PostgreSQL 9.5
    - Default cost parameter settings
    - **SeqScan** & **HashJoin**
        - **enable_seqscan = on, enable_hashjoin = on**
          and disables other methods
    - **IndexScan** & **NestLoopJoin**
        - **enable_indexscan = on, enable_nestloop = on**
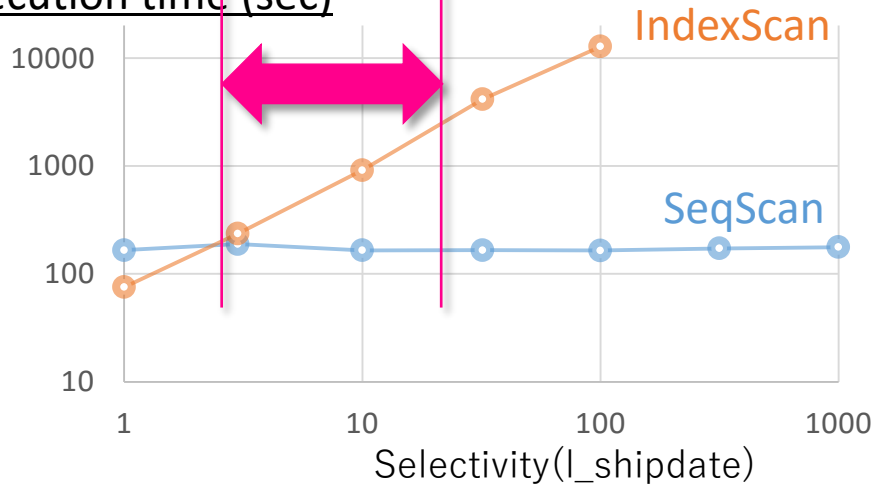          and disables other methods

# TPC-H Q.1: The Simplest Case

```
SELECT count(*), ... FROM lineitem
   WHERE l_shipdate BETWEEN [X] AND [Y]
```
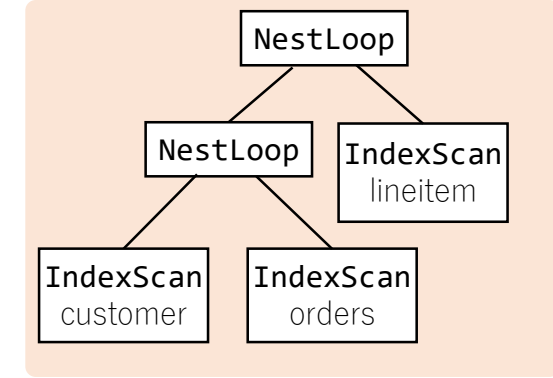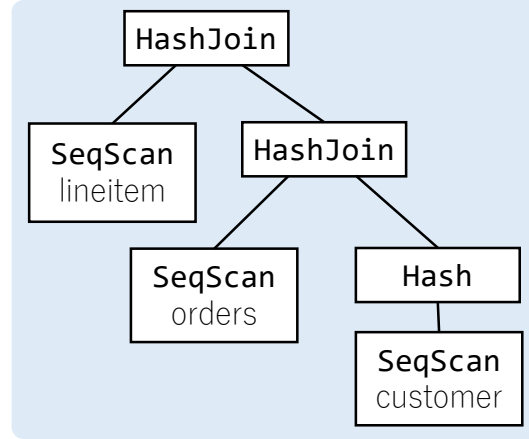
**Estimated cost**



Selectivity (l_shipdate)

**Execution time (sec)**
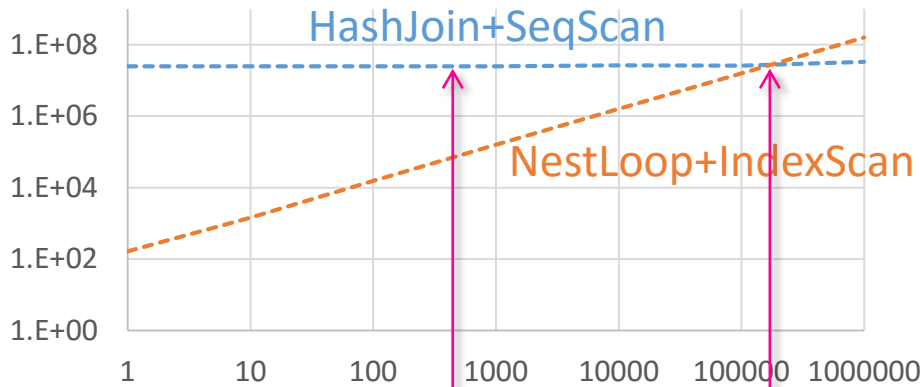


Selectivity (l_shipdate)

- Good trend estimation for each method

- Estimated break-event point is errorneus
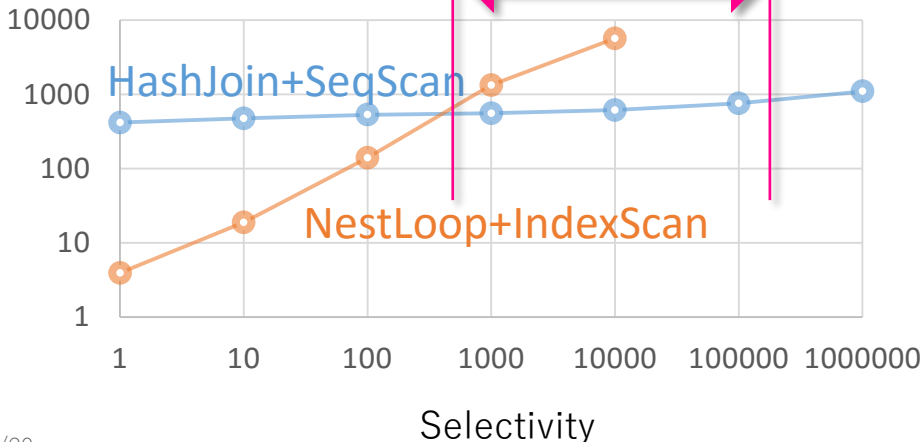  - IndexScan should be more expensive (need parameter calibration)

# TPC-H Q.3

```
        HashJoin
       /        \
SeqScan      HashJoin
lineitem    /        \
      SeqScan      Hash
      orders        |
                 SeqScan
                 customer
```

```
       NestLoop
      /        \
NestLoop    IndexScan
/     \     lineitem
IndexScan  IndexScan
customer   orders
```

## Estimated cost



HashJoin+SeqScan

NestLoop+IndexScan

Selectivity

| 1.E+08 |
| 1.E+06 |
| 1.E+04 |
| 1.E+02 |
| 1.E+00 |

1   10   100   1000   10000   100000   1000000

```
SELECT count(*), ...
   FROM customer, orders, lineitem
   WHERE c_custkey = o_custkey AND
         o_orderkey = l_orderkey AND
         c_custkey < [X] AND
         c_mktsegment = 'MACHINERY';
```
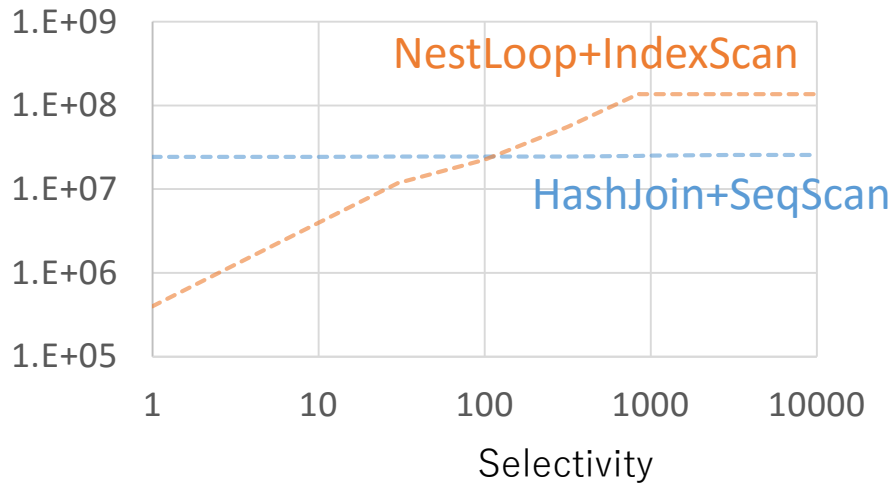
## Execution time (sec)



HashJoin+SeqScan

NestLoop+IndexScan

| 10000 |
| 1000 |
| 100 |
| 10 |
| 1 |

1   10   100   1000   10000   100000   1000000

Selectivity

Similar result as in Q.1

- Good trend estimation for each

- Erroneous break-event point without parameter calibration

# More Complex Case
# TPC-H Q.4: Semi-Join Query

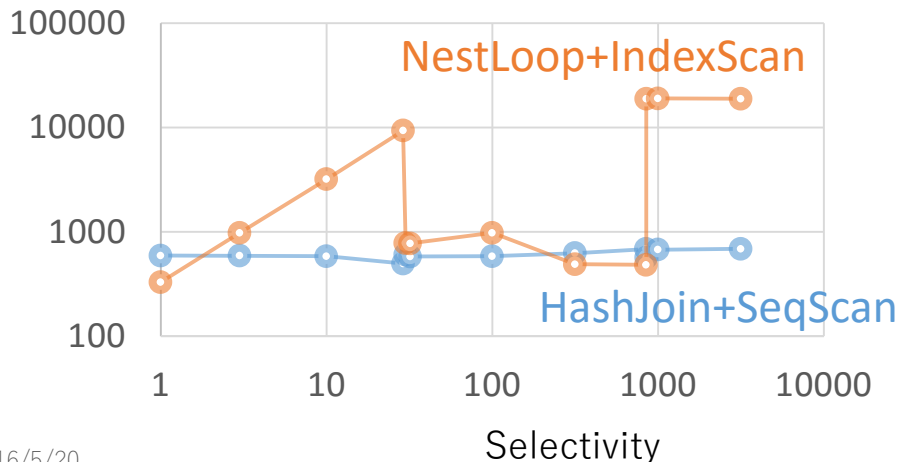### Estimated cost



```
SELECT count(*), ...
  FROM orders
  WHERE
    o_orderdate >= '1995-01-01' AND
    o_orderdate < '1995-01-01'
                  + interval '3 month' AND
    EXISTS(
      SELECT * FROM lineitem
       WHERE l_orderkey = o_orderkey
         AND l_commitdate < l_receiptdate)
```
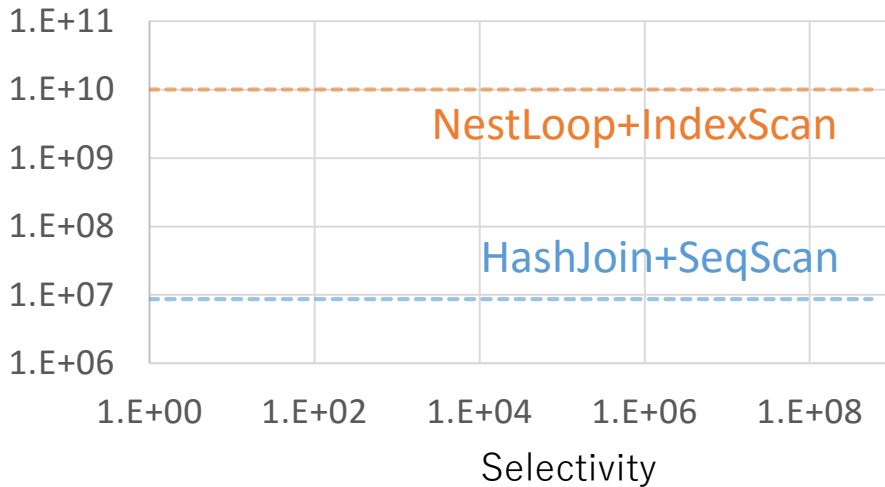
### Execution time (sec)



- Plan selection for semi-join tend to be unstable

# More Complex Case
# TPC-H Q.22: Anti-Join Query

**Estimated cost**



Selectivity

**Execution time (sec)**



Selectivity

```
SELECT count(*), ...
  FROM supplier, lineitem l1, orders, nation
  WHERE s_suppkey = l1.l_suppkey AND
        o_orderkey = l1.l_orderkey AND
        o_orderstatus = 'F' AND
        l1.l_receiptdate > l1.l_commitdate AND
    EXISTS (
      SELECT * FROM lineitem l2
        WHERE l2.l_orderkey = l1.l_orderkey
          AND l2.l_suppkey <> l1.l_suppkey)
  AND NOT EXIST (
      SELECT * FROM lineitem l3
        WHERE l3.l_orderkey = l1.l_orderkey
          AND l3.l_suppkey <> l1.l_suppkey
          AND l3.l_receiptdate > l3.l_commitdate)
  AND s_nationkey = n_nationkey
  AND n_name = 'JAPAN'
```

- Difficulties in overall cost trend estimation
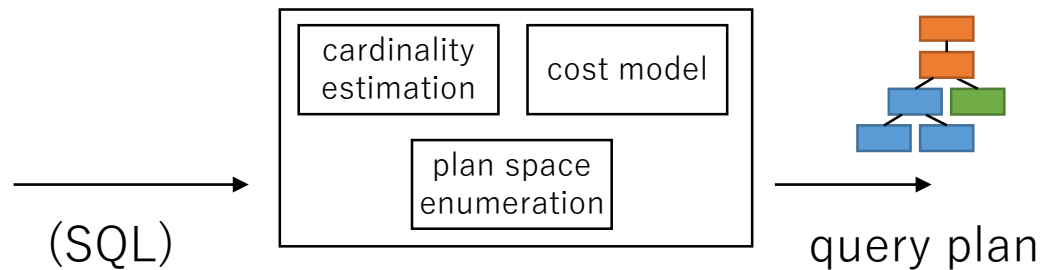
# Summary: PostgreSQL Optimizer

- Detailed look at cost modeling of basic methods
  - SeqScan, IndexScan
  - HashJoin, NestedLoopJoin
- Observation with TPC-H benchmark
  - Good cost trend estimation for simple join queries
    - Erroneous cheapest plan selection without parameter tuning
  - Difficulties with semi-join and anti-join queries

# Outline

- Introduction
- Theory: Query Optimization Framework
- Code: PostgreSQL Optimizer
- **Theory: Cutting-Edge Technologies Overview**
- **Summary**

# Cutting-Edge Technologies

- Traditional optimization was a "closed" problem



- "Rethink the contract" — *Surajit Chaudhuri*
  - Feedback from previous execution
  - Dynamic integration with execution

# Mid-query Re-optimization

[N. Kabra et.al., SIGMOD'98]

- Detects sub-optimality of executing query plan
  - Query plans are annotated for later estimation improvement
  - Runtime statistics collection
    - Statistics collector probes are inserted into operators of executing query plan
- Plan modification strategy
  - Discard current execution and re-optimize whole plan
  - Re-optimizer only subtree of the plan that are not started yet
  - Save partial execution result and generate new SQL using the result

# Plan Bouquet

- Generate a set of plans for each selectivity range
- Estimation improvement with runtime statistics collection

- Evaluation with PostgreSQL

# Bounding Impact of Estimation Error

**[T. Neumann et.al., BTW Conf '13]**

- "Uncertainty" analysis of cost estimation
  - Optimality sensitivity to estimation error

- Execute partially to reduce uncertainty

# Outline

- Introduction
- Theory: Query Optimization Framework
- Code: PostgreSQL Optimizer
- Theory: Cutting-Edge Technologies Overview
- Summary

# Summary

- Cost-based optimization framework
    - System-R style bottom-up optimization
    - Volcano style top-down optimization
- Detailed look at PostgreSQL optimizer
    - Cost modeling of basic scan and join method
    - Experiment with TPC-H benchmark
- Brief overview of cutting-edge technologies