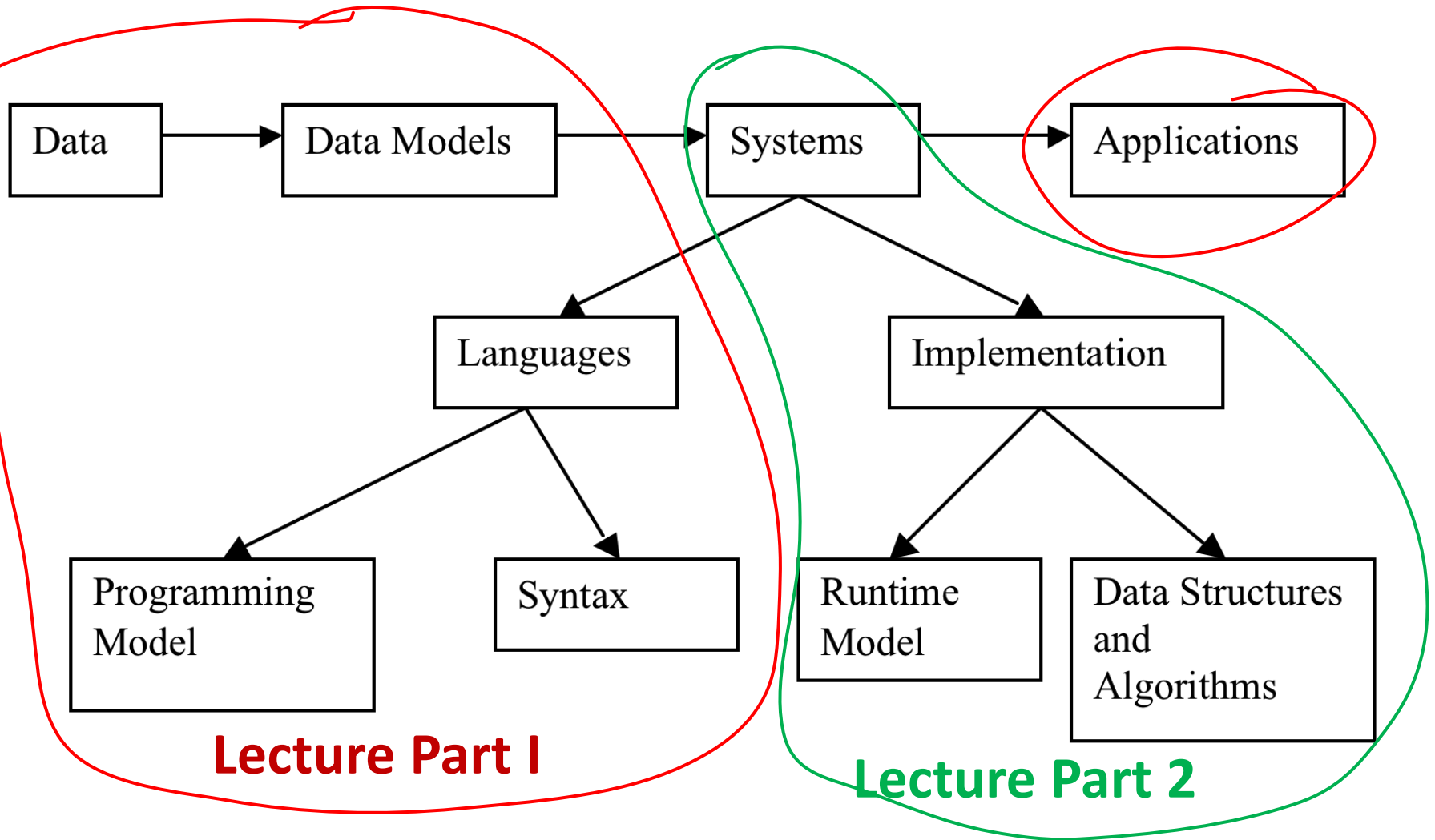# Module 4

# Implementation of XQuery

Part 0: Background on relational query processing

# The Data Management Universe



```
Data → Data Models → Systems → Applications

Systems → Languages
Systems → Implementation

Languages → Programming Model
Languages → Syntax

Implementation → Runtime Model
Implementation → Data Structures and Algorithms
```

**Lecture Part I**

**Lecture Part 2**

2

# What does a Database System do?

- Input: SQL statement
- Output: {tuples}

1. *Translate SQL into get/put requests to backend storage*
2. *Extract, process, transform tuples from blocks*

- Tons of optimizations
  - Efficient algorithms for SQL operators (hashing, sorting)
  - Layout of data on backend storage (clustering, free space)
  - Ordering of operators (small intermediate results)
  - Semantic rewritings of queries
  - Buffer management and caching
  - Parallel execution and concurrency
  - Outsmart the OS
  - Partitioning and Replication in distributed system
  - Indexing and Materialization
  - Load and admission control

- + Security + Durability + Concurrency Control + Tools

# XQuery: a mix of paradigms

- Query languages (~SQL)
- Functional programming languages  (~Haskell)
- Object-oriented query languages  (~OQL)
- Procedural languages  (~Java)
- Some new features : context sensitive semantics

- Processing XQuery involves
  - stealing from <u>all</u> other languages
  - plus specific innovations

# XQuery processing: old and new

## Functional programming

- ☑ Environment for expressions
- ☑ Expressions nested with full generality
- ☑ Lazy evaluation
- ☒ Data Model, schemas, type system, and query language
- ☒ Contextual semantics for expressions
- ☒ Side effects
- ☒ Non-determinism in logic operations, others
- ☒ *Streaming execution*
- ☒ *Logical/physical data mismatch, appropriate optimizations*

## Relational query (SQL)

- ☑ High level construct (FLWOR/Select-From-Where)
- ☑ Streaming execution
- ☑ Logical/physical data mismatch and the appropriate optimizations
- ☒ *Data Model, schemas, type system, and query language*
- ☒ *Expressive power*
- ☒ Error handling
- ☒ 2 valued logic

# XQuery processing: old and new

**Object-oriented query languages (OQL)**

- ☑ Expressions nested with full generality
- ☑ Nodes with node/object identity
- ☒ Topological order for nodes
- ☒ *Data Model, schemas, type system, and query language*
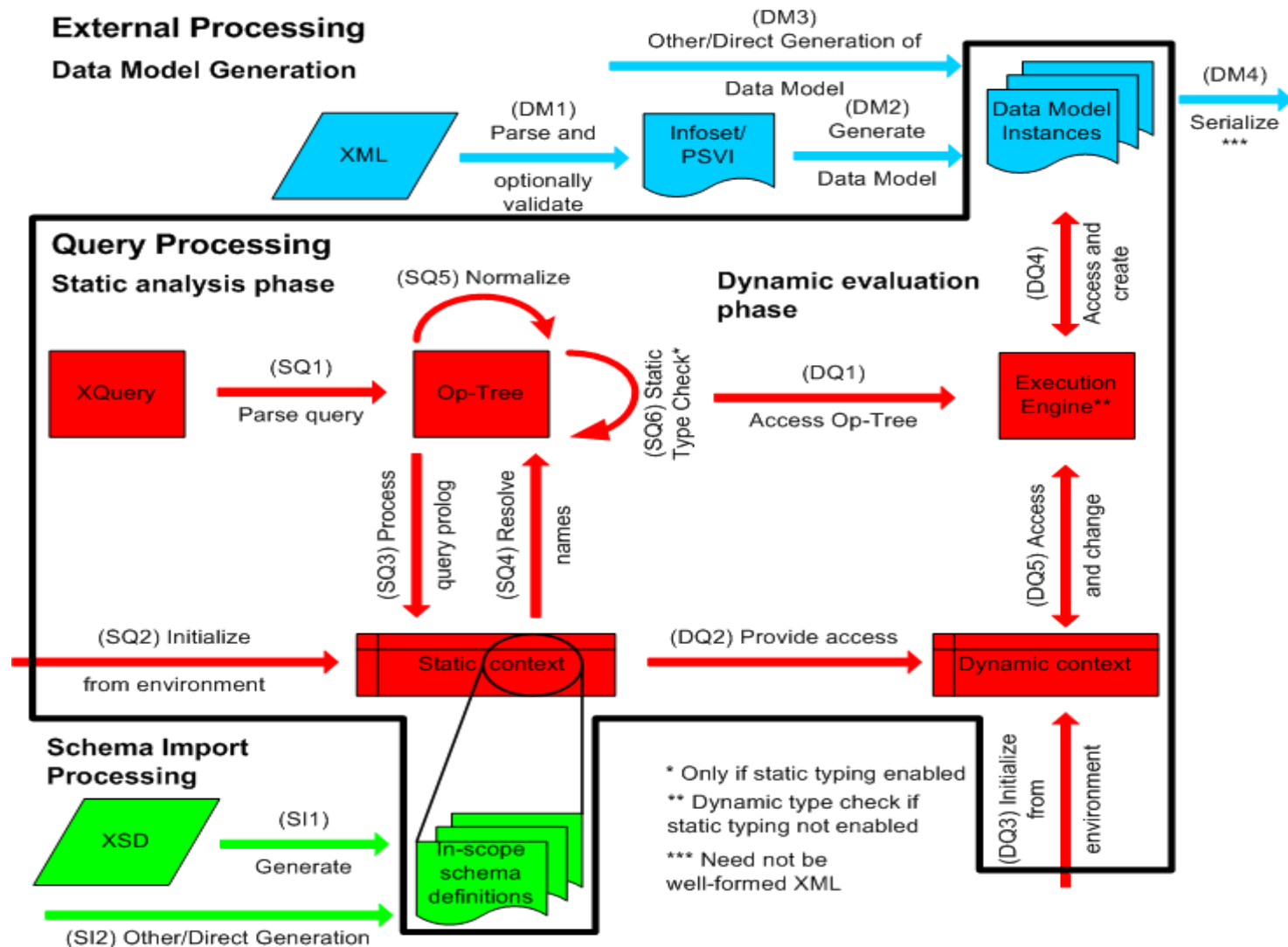- ☒ *Side effects*
- ☒ *Streaming execution*

**Imperative languages (e.g. Java)**

- ☑ Side effects
- ☑ Error handling
- ☒ *Data Model, schemas, type system, and query language*
- ☒ Non-determinism for logic operators
- ☒ *Lazy evaluation and streaming*
- ☒ *Logical/physical data mismatch and the appropriate optimizations*
- ☒ *Possibility of handling large volumes of data*
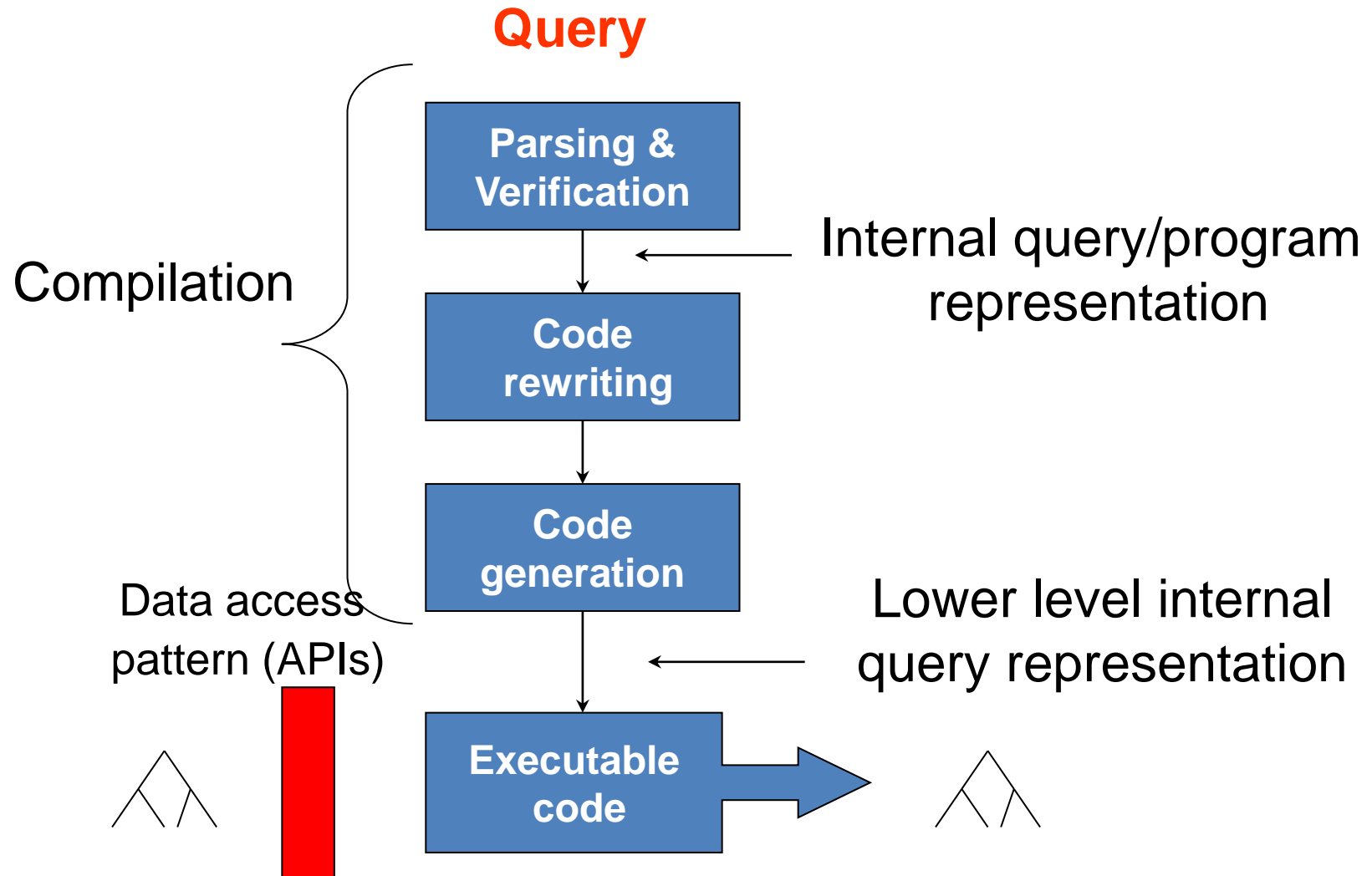
# Aspects of XQuery Implementation

- **Compile Time + Optimizations**
  - Operator Models
  - Query Rewrite
  - Runtime + Query Execution
- **XML Data Representation**
  - XML Storage
  - XML Indexes
  - Compression + Binary XML

# XQuery Processing Model

# Architecture of (X)Query Processor

**Query**

Parsing & Verification

Compilation

Code rewriting

Internal query/program representation

Code generation

Data access pattern (APIs)

Lower level internal query representation

Executable code

# Backgrounds from the database world

- Database management systems provides
  - a success story for building large-data, declarative infrastructures
  - Blueprints on architecture and algorithms
- No class at Uni Freiburg (explicitly) teaches these contents (as opposed to e.g., compiler construction etc)
- Quick tour through relational concepts:
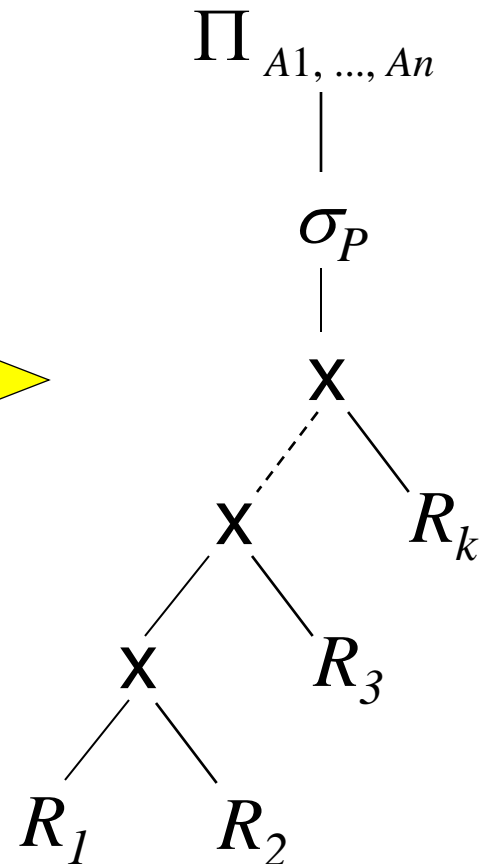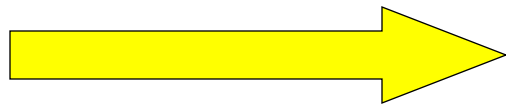  - Algebra
  - Query Processing
  - Optimization

# SQL -> Relational Algebra

SQL

Relational Algebra

$$\Pi_{A1, \dots, An}(\sigma_P (R_1 \times \dots \times R_k))$$

$$\Pi_{A1, \dots, An}$$

| 

$$\sigma_P$$

|

$\times$

**select** $A_1, \dots, A_n$

**from** $R_1, \dots, R_k$

**where** $P$;

$\times$    $R_k$
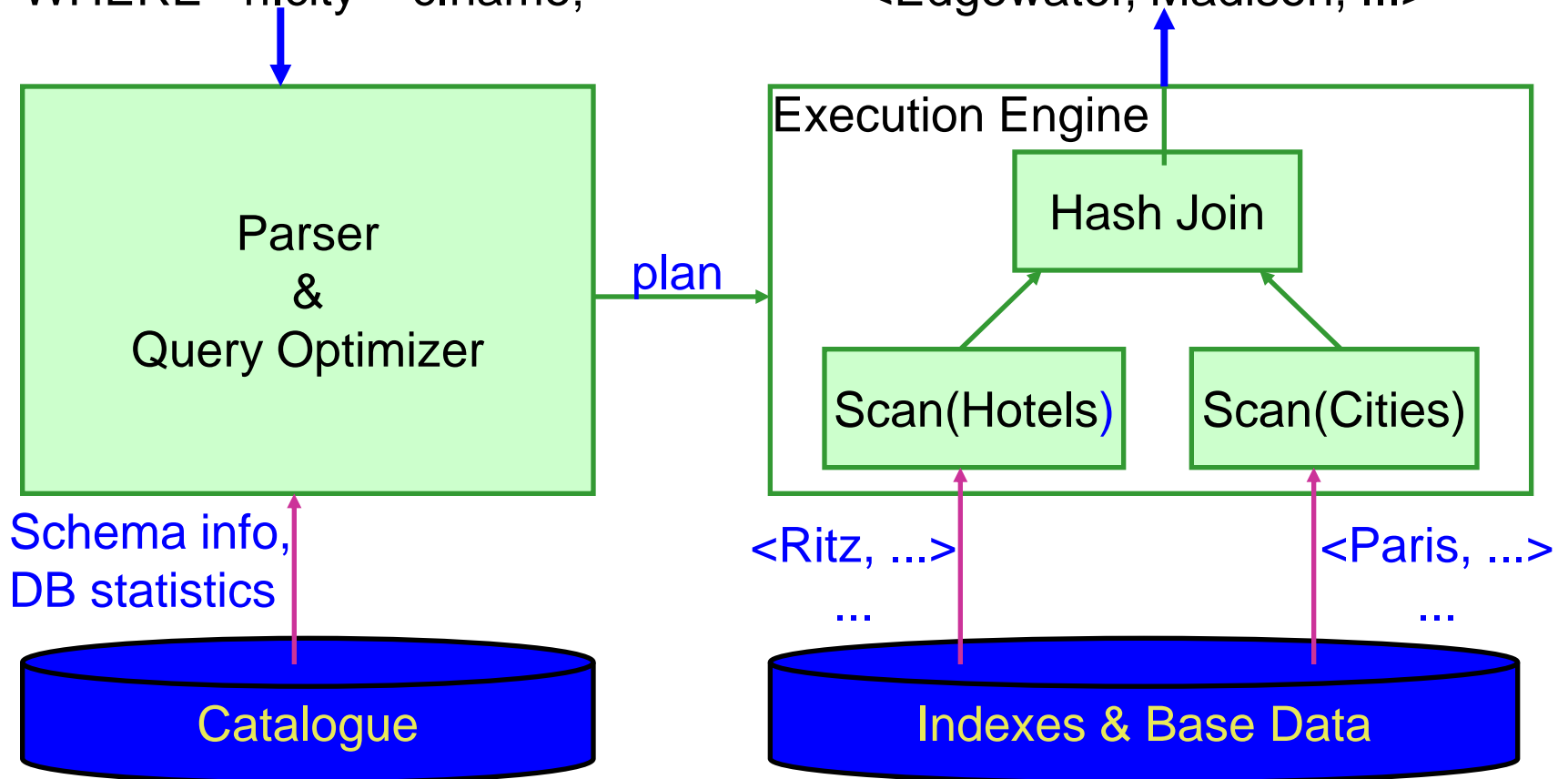
$\times$    $R_3$

$R_1$    $R_2$

# Algorithms for Rel. Algebra

- Table Access
  - scan     (load each page at a time)
  - index scan     (if index available)
- Sorting
  - Two-phase external sorting
- Joins
  - (Block) nested-loops
  - Index nested-loops
  - Sort-Merge
  - Hashing (many variants)
- Group-by  (~ self-join)
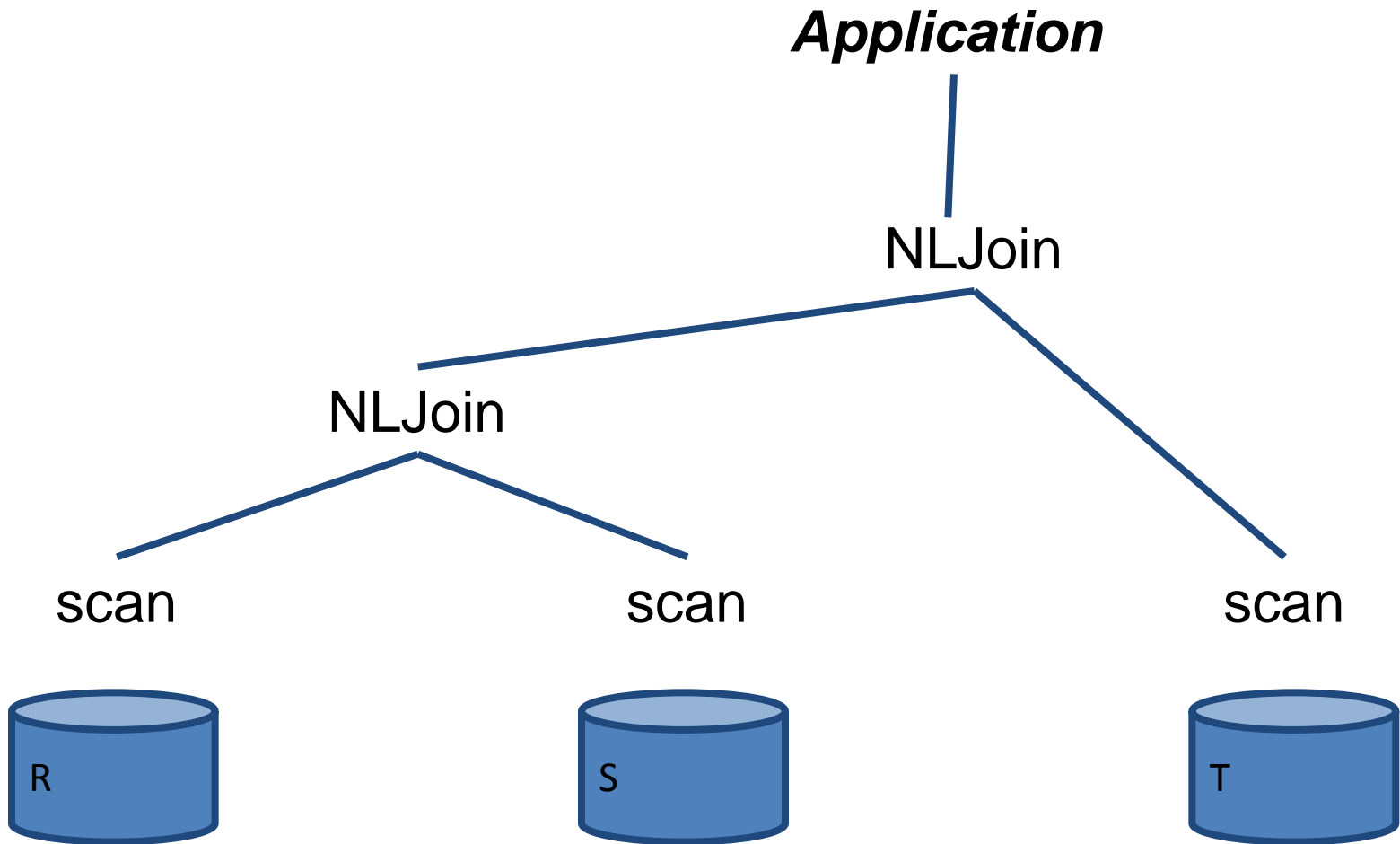  - Sorting, Hashing

# SQL Query Processing 101

```
SELECT    *
FROM      Hotels h, Cities c
WHERE     h.city = c.name;
```
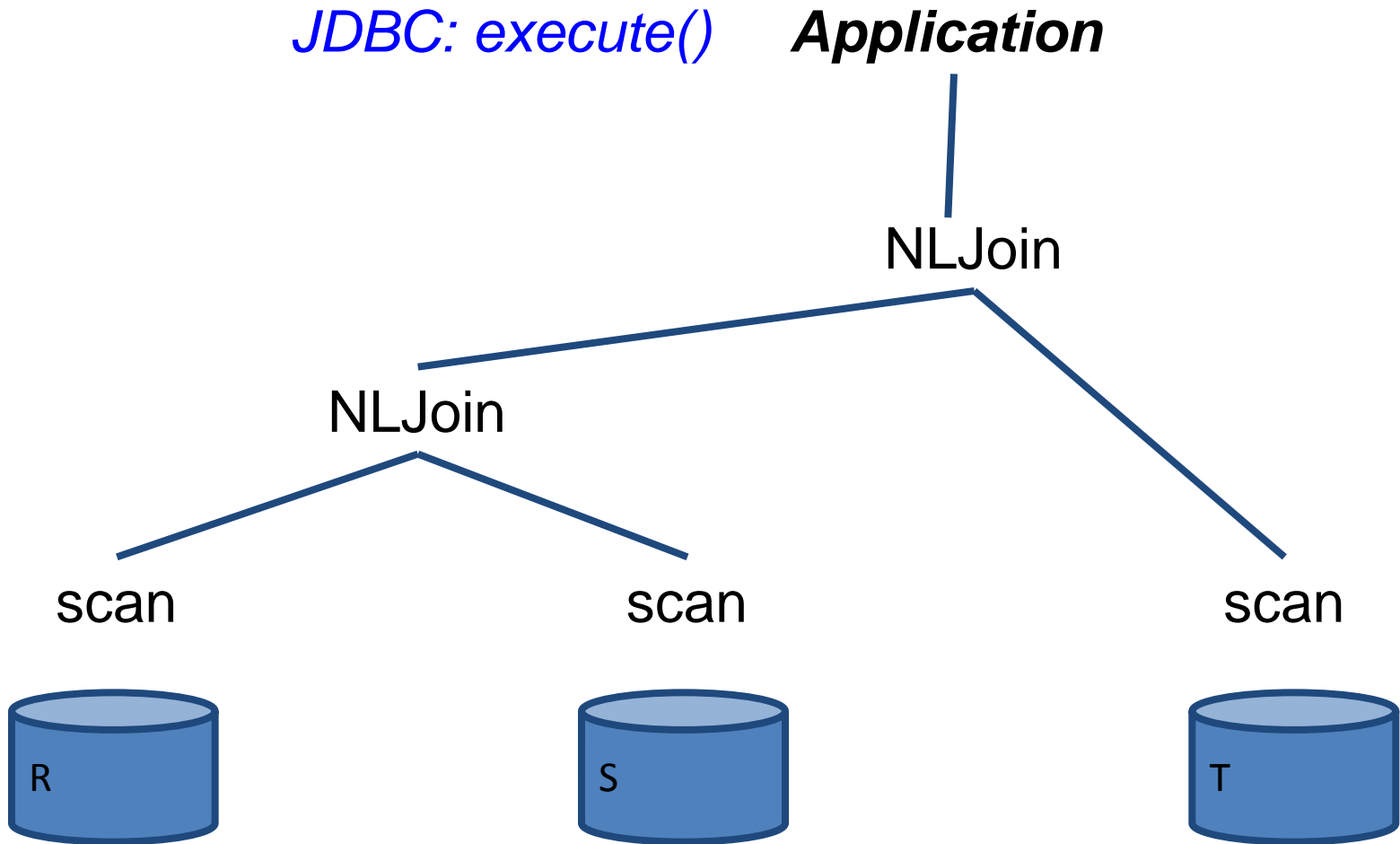
<Ritz, Paris, ...>
<Baur au Lac, Zurich, ...>
<Edgewater, Madison, ...>

Parser
&
Query Optimizer

plan

Execution Engine

Hash Join

Scan(Hotels)          Scan(Cities)

Schema info,
DB statistics

<Ritz, ...>                          <Paris, ...>
...                                  ...

Catalogue

Indexes & Base Data

# Iterator Model

- Plan contains many operators
  - Implement each operator indepently
  - Define generic interface for each operator
  - Each operator implemented by an *iterator*
- Three methods implemented by each iterator
  - open(): initialize the internal state (e.g., buffer)
  - char* next(): produce the next result tuple
  - close(): clean-up (e.g., release buffer)
- N.B. Modern DBMS use a Vector Model
  - next() returns a set of tuples
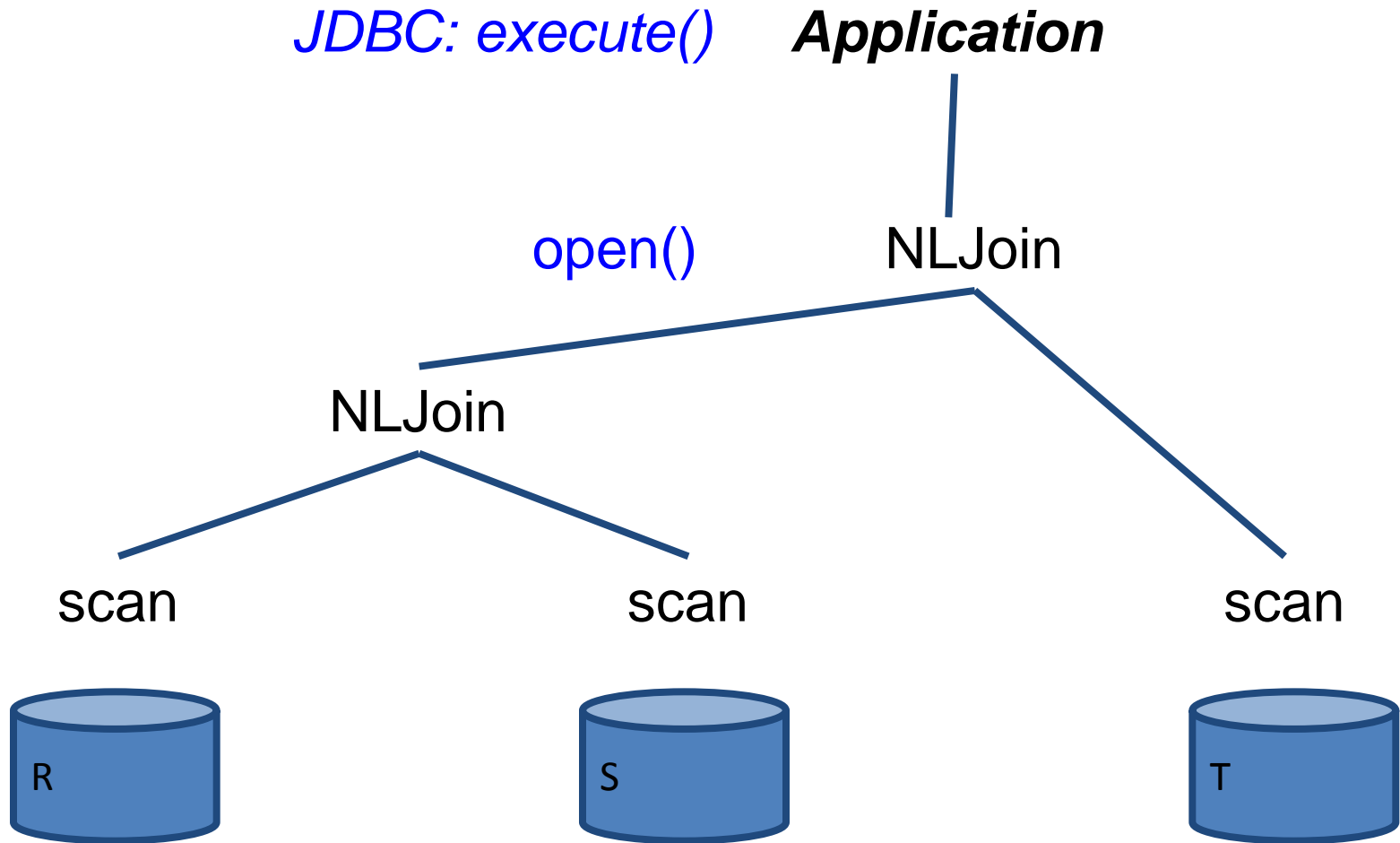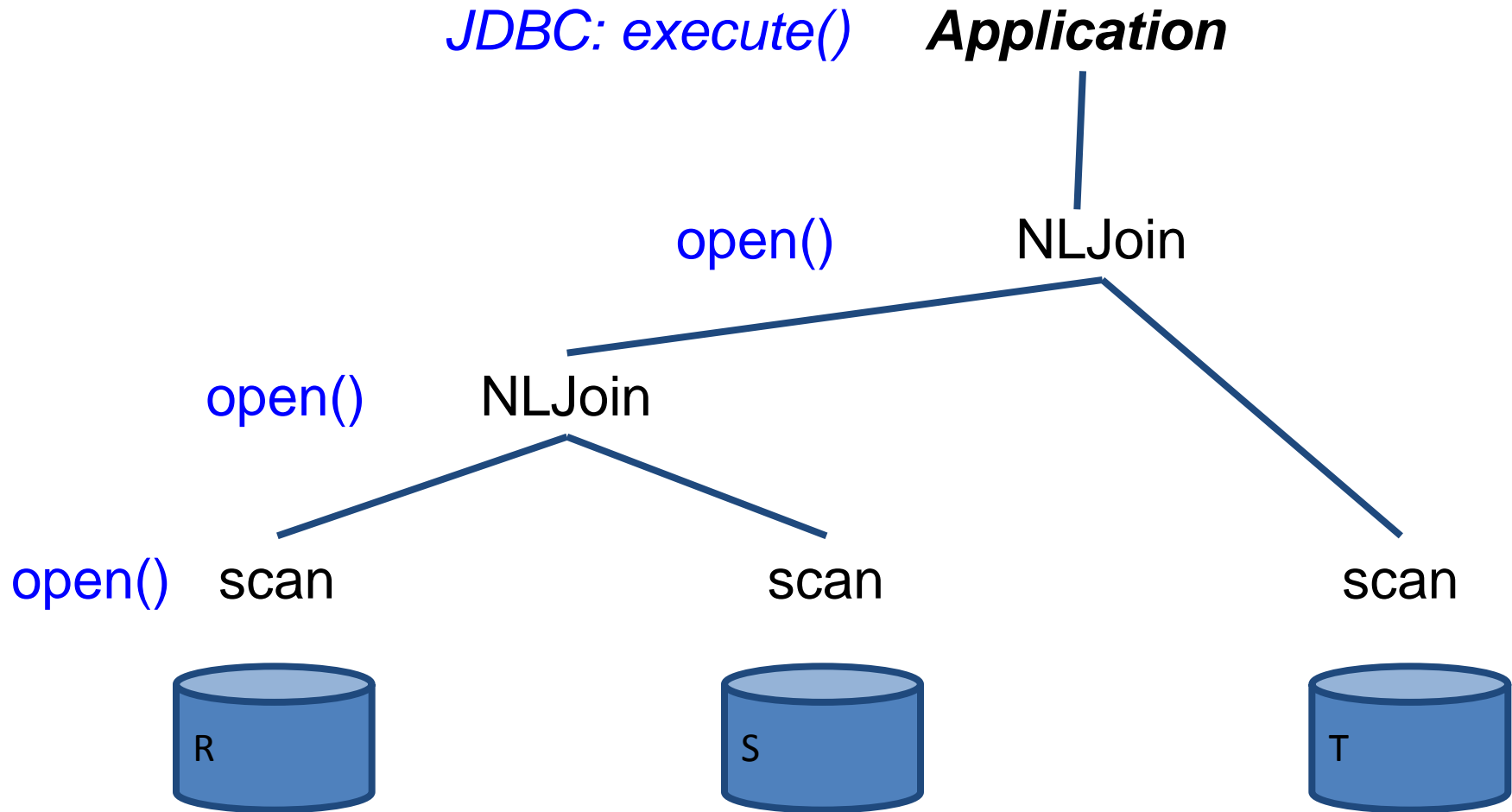  - Why is that better?

# Iterator Model at Work

*Application*

NLJoin
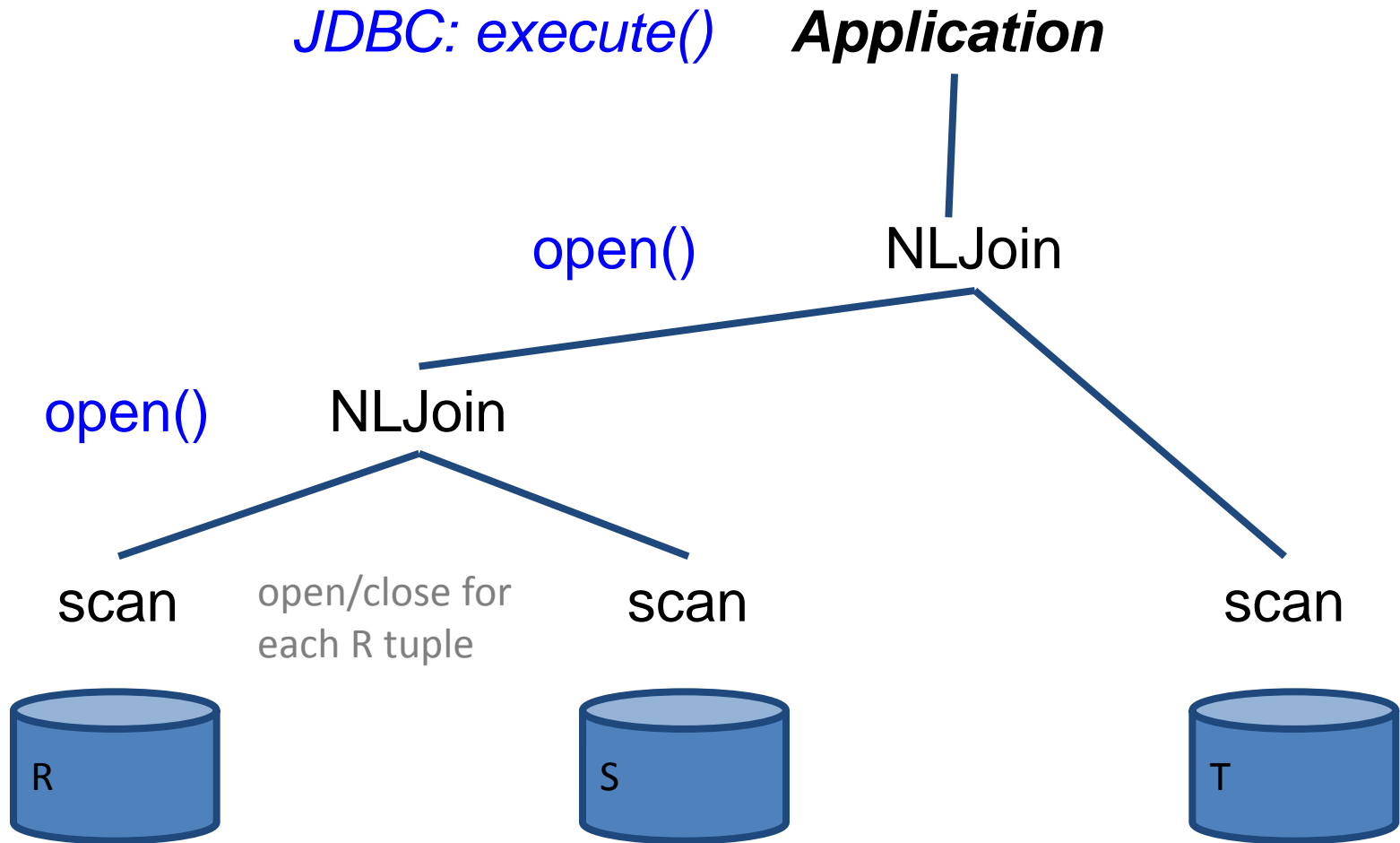
NLJoin

scan                    scan                    scan

R                       S                       T

# Iterator Model at Work

*JDBC: execute()*  **Application**

NLJoin

NLJoin

scan          scan          scan

R          S          T

# Iterator Model at Work

*JDBC: execute()*     **Application**

open()     NLJoin

NLJoin

scan        scan        scan

R       S       T

# Iterator Model at Work

*JDBC: execute()*   **Application**

open()   NLJoin

open()   NLJoin

open()   scan          scan                    scan

R              S                    T

# Iterator Model at Work

*JDBC: execute()*    ***Application***

open()    NLJoin

open()    NLJoin

scan    open/close for each R tuple    scan      scan

R       S       T

# Iterator Model at Work

*JDBC: execute()*    **Application**

open()    NLJoin

NLJoin

scan        scan    open/close for each R,S tuple    scan

R          S          T

# Iterator Model at Work

*JDBC: next()*     ***Application***

next()     NLJoin

next()     NLJoin

next()    scan      scan      scan

R      S      T

# Iterator Model at Work

*JDBC: next()*     **Application**

next()                    NLJoin

next()    NLJoin

r1    scan          scan          scan

R            S            T

# Iterator Model at Work

*JDBC: next()*      **Application**

next()      NLJoin

next()      NLJoin

r1

scan      open()      scan              scan

R                    S                    T

# Iterator Model at Work

*JDBC: next()*    **Application**

next()    NLJoin

next()    NLJoin

r1

scan    next()    scan    scan

R    S    T

# Iterator Model at Work

*JDBC: next()*   **Application**

next()   NLJoin

next()   NLJoin

r1

scan   s1   scan   scan

R   S   T

# Iterator Model at Work

*JDBC: next()*          **Application**

next()          NLJoin

next()     NLJoin

r1

scan          next()     scan          scan

R          S          T

# Iterator Model at Work



*JDBC: next()*     **Application**

next()     NLJoin

next()     NLJoin

r1

scan     s2     scan     scan

R     S     T

# Iterator Model at Work

*JDBC: next()*       **Application**

next()       NLJoin

r1, s2    NLJoin

scan                scan                scan

R                   S                   T

# Iterator Model at Work

*JDBC: next()*    **Application**

next()    NLJoin

r1, s2

NLJoin

scan    scan    open()    scan

R    S    T

# Iterator Model at Work

*JDBC: next()*     **Application**

next()     NLJoin

r1, s2

NLJoin

scan          scan     next()     scan

R                S                    T

# Iterator Model at Work

*JDBC: next()*     **Application**

next()     NLJoin

r1, s2

NLJoin

scan     scan     t1     scan

R     S     T

# Iterator Model at Work



*JDBC: next()*   **Application**

r1, s2, t1   NLJoin

NLJoin

scan   scan   scan

R   S   T

… r2, r3, …   … s3, s4, …   … t2, t3, …

# Iterators: Easy & Costly

- Principle
  - data flows bottom up in a plan (i.e. operator tree)
  - control flows top down in a plan
- Advantages
  - generic interface for all operators: great information hiding
  - easy to implement iterators (clear what to do in any phase)
  - works well with JDBC and embedded SQL
  - supports DBmin and other buffer management strategies
  - no overheads in terms of main memory
  - supports pipelining: great if only subset of results consumed
  - supports parallelism and distribution: add special iterators
- Disadvantages
  - high overhead of method calls
  - poor instruction cache locality

# Compiler: Optimizations

- Goals:
  1. Reduce the *level of abstraction*
  2. Reduce the *execution cost*
- Concepts
  - Code representation (e.g., algebras)
  - Code transformations (e.g., rules)
  - Cost transformation policy (e.g., enumeration)
  - Code cost estimation

# SQL -> Relational Algebra

SQL

Relational Algebra

$$\Pi_{A1,\,...,\,An}(\sigma_P\,(R_1 \times ... \times R_k))$$

$\Pi_{A1,\,...,\,An}$

|

$\sigma_P$

|

$\times$

$\times$      $R_k$

$\times$      $R_3$

$\times$

$R_1$    $R_2$

**select** $A_1, ..., A_n$

**from** $R_1, ..., R_k$

**where** $P;$

# SQL -> QGM

SQL

QGM

**select** *a*

**from** *R*

**where** *a in (***select** *b*

    **from** *S);*



$\Pi_a$

$\sigma$

*R*

in

$\Pi_b$

*S*

# Parser

- Generates rel. alg. tree for each sub-query
  - constructs graph of trees: Query Graph Model nodes are subqueries
  - edges represent relationships between subqueries
- Extended rel. algebra because SQL more than RA
  - GROUP BY: $\Gamma$ operator
  - ORDER BY: sort operator
  - DISTINCT: can be implemented with $\Gamma$ operator

# SQL -> Relational Algebra

SQL

Relational Algebra

$$\Pi_{A1,\,...,\,An}(\sigma_P (R_1 \times ... \times R_k))$$

$$\Pi_{A1,\,...,\,An}$$

$$\sigma_P$$

**select** $A_1, ..., A_n$

**from** $R_1, ..., R_k$

**where** $P;$
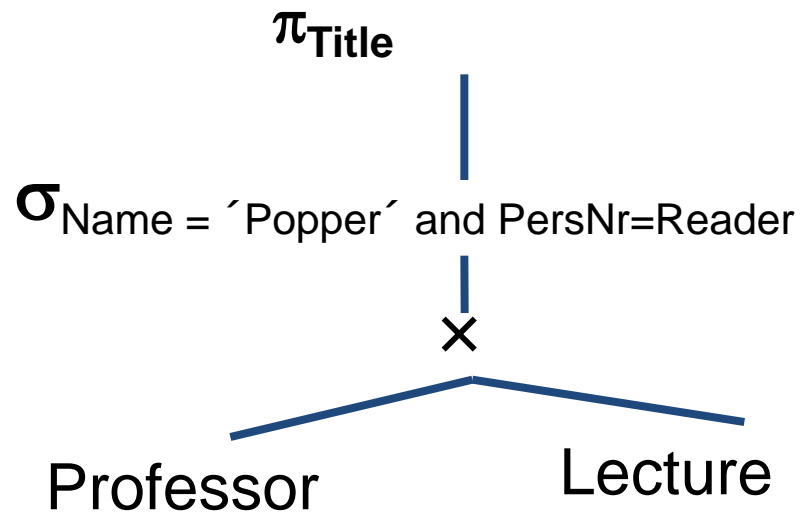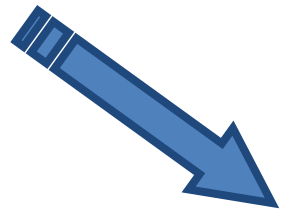
$\times$

$\times$    $R_k$

$\times$    $R_3$
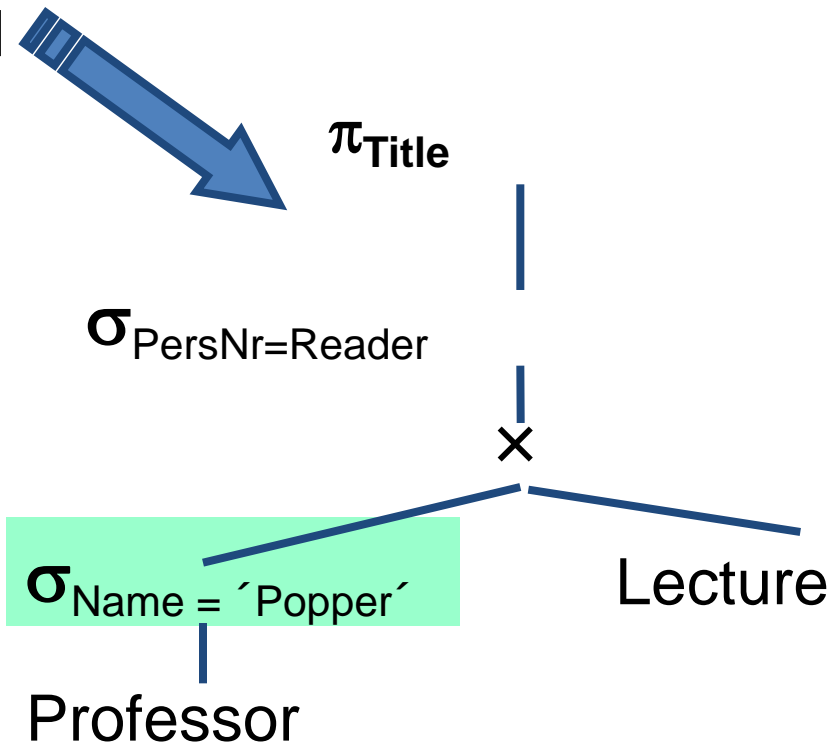
$R_1$    $R_2$

# Example: SQL -> Relational Algebra

**select** Title

**from** Professor, Lecture

**where** Name = ´Popper´ **and**

PersNr = Reader

$\pi_{\textbf{Title}}$

$\sigma_{\text{Name = ´Popper´ and PersNr=Reader}}$

$\times$

Professor              Lecture

$\pi_{\textbf{Title}} (\sigma_{\text{Name = ´Popper´ and PersNr=Reader}} (\text{Professor} \times \text{Lecture}))$

# First Optimization: Push-down $\sigma$

**select** Title

**from** Professor, Lecture
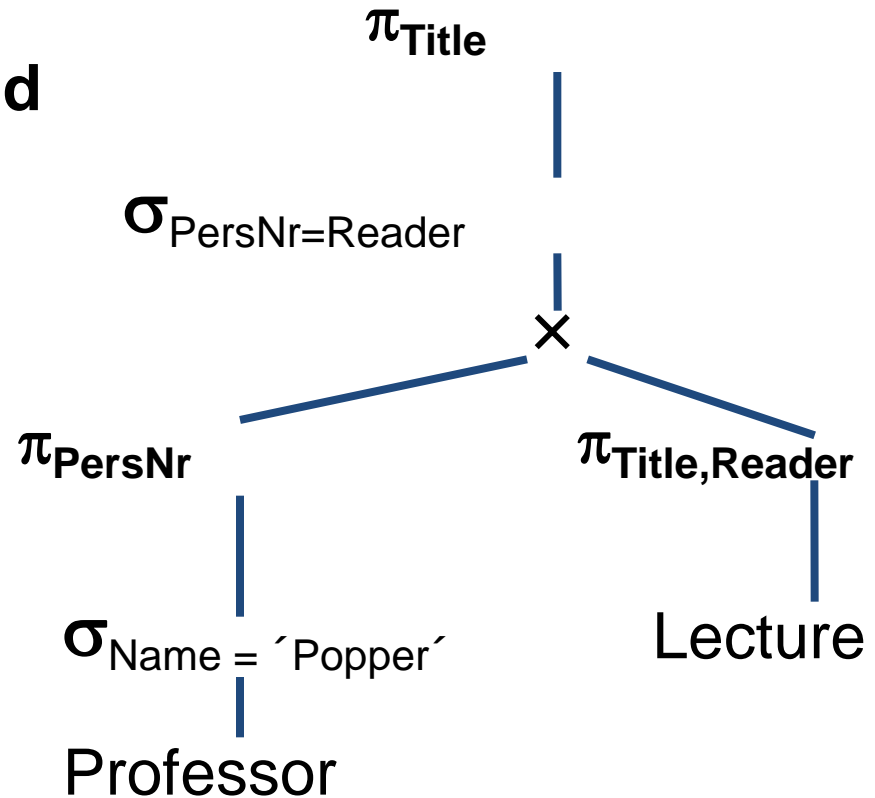
**where** Name = ´Popper´ **and**

PersNr = Reader

$\pi_{\textbf{Title}}$

$\sigma_{\text{PersNr=Reader}}$

$\times$

$\sigma_{\text{Name = ´Popper´}}$

Lecture

Professor
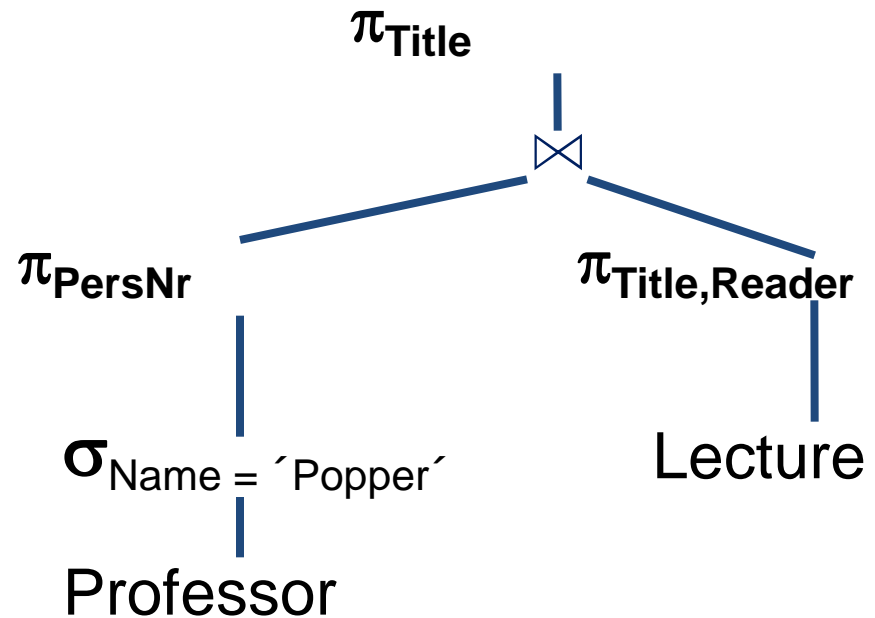
$\pi_{\textbf{Title}} (\sigma_{\text{PersNr=Reader}} ((\sigma_{\text{Name = ´Popper´}} \text{ Professor}) \times \text{Lecture}))$

# Push-down $\pi$

**select** Title
**from** Professor, Lecture
**where** Name = ´Popper´ **and**
$\qquad$ PersNr = Reader

$\pi_{\textbf{Title}}$

$\sigma_{\text{PersNr=Reader}}$

$\times$

$\pi_{\textbf{PersNr}}$

$\pi_{\textbf{Title,Reader}}$

$\sigma_{\text{Name = ´Popper´}}$

Lecture

Professor

# Correctness: Push-down $\pi$

- $\pi_{\textbf{Title}} \left( \sigma_{\text{PersNr=Reader}} \left( \left( \sigma_{\text{Name = ´Popper´}} \text{Professor} \right) \times \text{Lecture} \right) \right)$

  (composition of projections)

- $\pi_{\textbf{Title}} \left( \pi_{\textcolor{red}{\textbf{Title,PersNr,Reader}}} \left( \sigma_{...} \left( \left( \sigma_{...} \text{Professor} \right) \times \text{Lecture} \right) \right) \right)$

  (commutativity of $\pi$ and $\sigma$)

- $\pi_{\textbf{Title}} \left( \sigma_{...} \left( \pi_{\textcolor{red}{\textbf{Title,PersNr,Reader}}} \left( \left( \sigma_{...} \text{Professor} \right) \times \text{Lecture} \right) \right) \right)$

  (commutativity of $\pi$ and $\sigma$)

- $\pi_{\textbf{Title}} \left( \sigma_{...} \left( \pi_{\textcolor{red}{\textbf{PersNr}}} \left( \sigma_{...} \text{Professor} \right) \times \pi_{\textcolor{red}{\textbf{Title,Reader}}} \left( \text{Lecture} \right) \right) \right)$

42

# Push down $\pi$

- Correctness (see previous slide – example generalizes)

- Why is it good?  ( almost same reason as for $\sigma$)
  - reduces size of intermediate results
  - but: only makes sense if results are materialized; e.g. sort
    - does not make sense if pointers are passed around in iterators

# Third Optimization: $\sigma$ + x = $\bowtie$

**select** Title

**from** Professor, Lecture

**where** Name = ´Popper´ **and**

      PersNr = Reader

$\pi_{\text{Title}}$

$\bowtie$

$\pi_{\text{PersNr}}$            $\pi_{\text{Title,Reader}}$

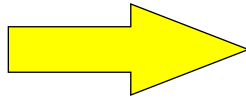$\sigma_{\text{Name = ´Popper´}}$       Lecture

Professor

# Third Optimization: $\sigma$ + x = $\bowtie$

- Correctness by definition of $\bowtie$ operator
- Why is this good?
  - x always done using nested-loops algorithm
    - $\bowtie$ can also be carried out using hashing, sorting, index support
    - choice of better algorithm may result in huge wins
  - x produces large intermediate results
    - results in a huge number of „next()" calls in iterator model
    - method calls are expensive
- Selection, projection push-down are no-brainers
  - make sense whenever applicable
  - do not need a cost model to decide how to apply them
  - (exception: expensive selections, projections with UDF)
  - done in a phase called query rewrite, based on rules
- More complex query rewrite rules…

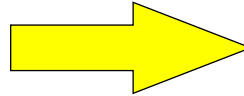# Unnesting of Views

- ## Example: Unnesting of Views

```
select A.x
from   A
where y in
   (select y from B)
```

→

```
select A.x
from   A, B
where A.y = B.y
```

- ## Example: Unnesting of Views

```
select A.x
from   A
where exists
   (select * from B where A.y = B-y)
```
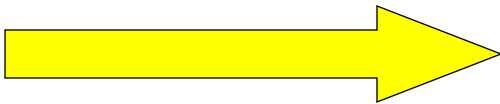
→

```
select A.x
from   A, B
where A.y = B.y
```

- ## Is this correct? Why is this better?
  - (not trivial at all!!!)

# Query Rewrite

- Example: Predicate Augmentation
  ```
  select *
  from   A, B, C
  where  A.x = B.x
   and B.x = C.x
  ```

  ```
  select *
  from   A, B, C
  where  A.x = B.x
   and B.x = C.x
  and A.x = C.x
  ```

**Why is that useful?**

# Pred. Augmentation: Why good?

A (odd numbers)

| ... | x |
|-----|-----|
| ... | 1 |
| ... | 3 |
| ... | 5 |
| ... | ... |

B (all numbers)

| ... | x |
|-----|-----|
| ... | 1 |
| ... | 2 |
| ... | 3 |
| ... | ... |

C (even numbers)

| ... | x |
|-----|-----|
| ... | 2 |
| ... | 4 |
| ... | 6 |
| ... | ... |

- Cost$((A \bowtie C) \bowtie B) <$ Cost$((A \bowtie B) \bowtie C)$
  - get second join for free
- Query Rewrite does not know that, …
  - but it knows that it might happen and hopes for optimizer…
- Codegen gets rid of unnecessary predicates (e.g., A.x = B.x)

# Query Optimization

- Two tasks
  - Determine order of operators
  - Determine algorithm for each operator (hash vs sort)
- Components of a query optimizer
  - Search space
  - Cost model
  - Enumeration algorithm
- Working principle
  - Enumerate alternative plans
  - Apply cost model to alternative plans
  - Select plan with lowest expected cost

# Query Opt.: Does it matter?

- A x B x C
  - size(A) = 10,000
  - size(B) = 100
  - size(C) = 1
  - cost(X x Y) = size(X) + size(Y)

- cost( (A x B) x C) = **1,010,001**
  - cost(A x B) = 10,100
  - cost(X x C) = 1,000,001    with X = A x B

- cost ( A x (B x C)) = **10,201**
  - cost(B x C) = 101
  - cost(A x X) = 10,100    with X = B x C

# Query Opt.: Does it matter?

- A x B x C
  - size(A) = 1000
  - size(B) = 1
  - size(C) = 1
  - **cost(X x Y) = size(X) * size(Y)**

- cost( (A x B) x C) = **2000**
  - cost(A x B) = 1000
  - cost(X x C) = 1000      with X = A x B

- cost ( A x (B x C)) = **1001**
  - cost(B x C) = 1
  - cost(A x X) = 1000      with X = B x C

# Search Space: Rel. Algebra

- Associativity of joins:

    $(A \bowtie B) \bowtie C = A \bowtie (B \bowtie C)$

- Commutativity of joins:

    $A \bowtie B = B \bowtie A$

- Many more rules
    - see Kemper/Eickler or Garcia-Molina text books

- What is better:  $A \bowtie B$ or $B \bowtie A$?
    - it depends
    - need cost model to make decision

# Search Space: Group Bys

SELECT … FROM R, S WHERE R.a = S.a GROUP BY R.a, S.b;

- $\Gamma_{R.a, S.b}(R \bowtie S)$
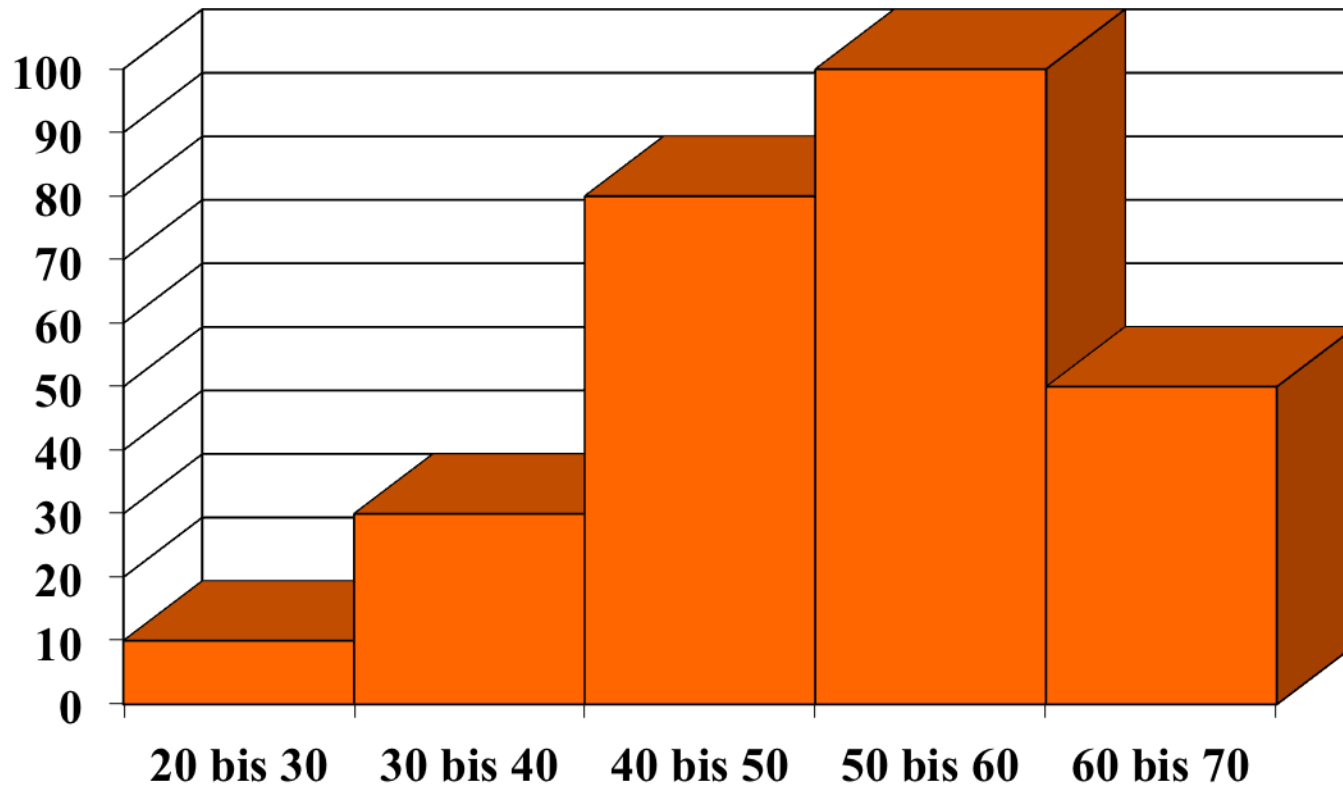
- $\Gamma_{S.b}(\Gamma_{R.a}(R) \bowtie S)$

- Often, many possible ways to split & move group-bys
  – again, need cost model to make right decisions

# Cost Model

- Cost Metrics
  - Response Time (consider parallelism)
  - Resource Consumption: CPU, IO, network
  - $  (often equivalent to resource consumption)
- Principle
  - Understand algorithm used by each operator (sort, hash, …)
    - estimate available main memory buffers
    - estimate the size of inputs, intermediate results
  - Combine cost of operators:
    - sum for resource consumption
    - max for response time (but keep track of bottlenecks)
- Uncertainties
  - estimates of buffers, interference with other operators
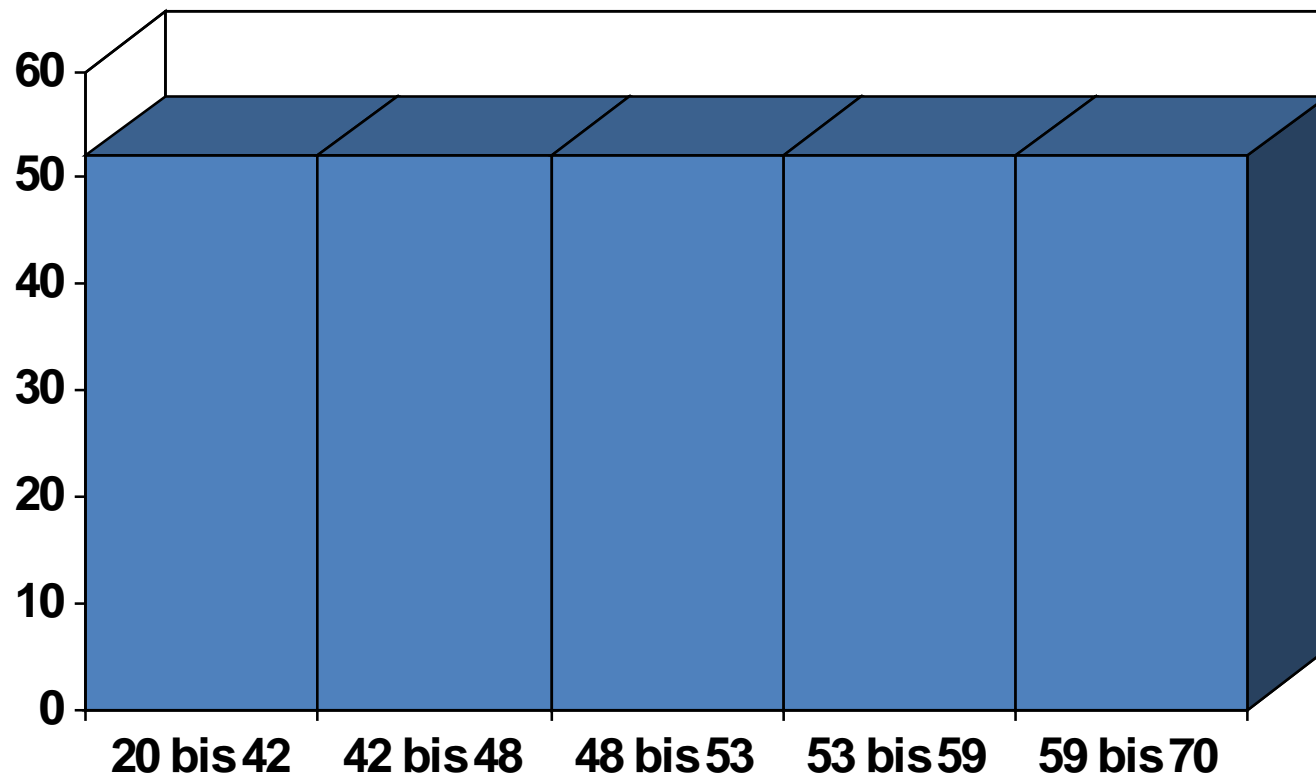  - estimates of intermediate result size (histograms)

# Equi-Width Histogram

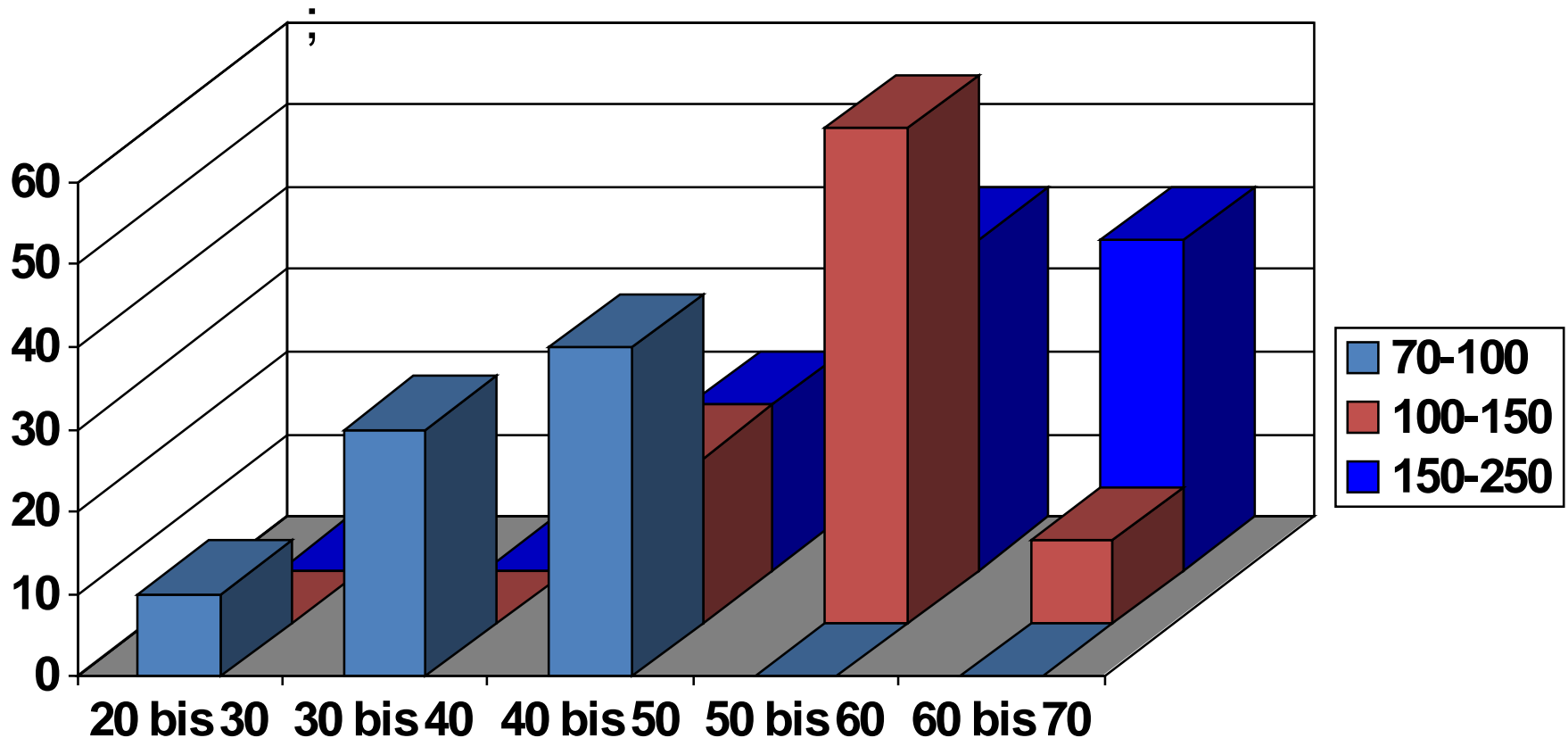SELECT * FROM person WHERE 25 < age < 40;

# Equi-Depth Histogram

SELECT * FROM person WHERE 25 < age < 40;

# Multi-Dimensional Histogram



SELECT * FROM person
WHERE 25 < age < 40 AND salary > 200;
;

# Enumeration Algorithms

- Query Optimization is NP hard
  - even ordering or Cartesian products is NP hard
  - in general impossible to predict complexity for given query
- Overview of Algorithms
  - Dynamic Programming (good plans, exp. complexity)
  - Greedy heuristics (e.g., highest selectivity join first)
  - Randomized Algorithms (iterative improvement, Sim.An., …)
  - Other heuristics (e.g., rely on hints by programmer)
  - Smaller search space (e.g., deep plans, limited group-bys)
- Products
  - Dynamic Programming used by many systems
  - Some systems also use greedy heuristics in addition