# CS 525: Advanced Database Organisation
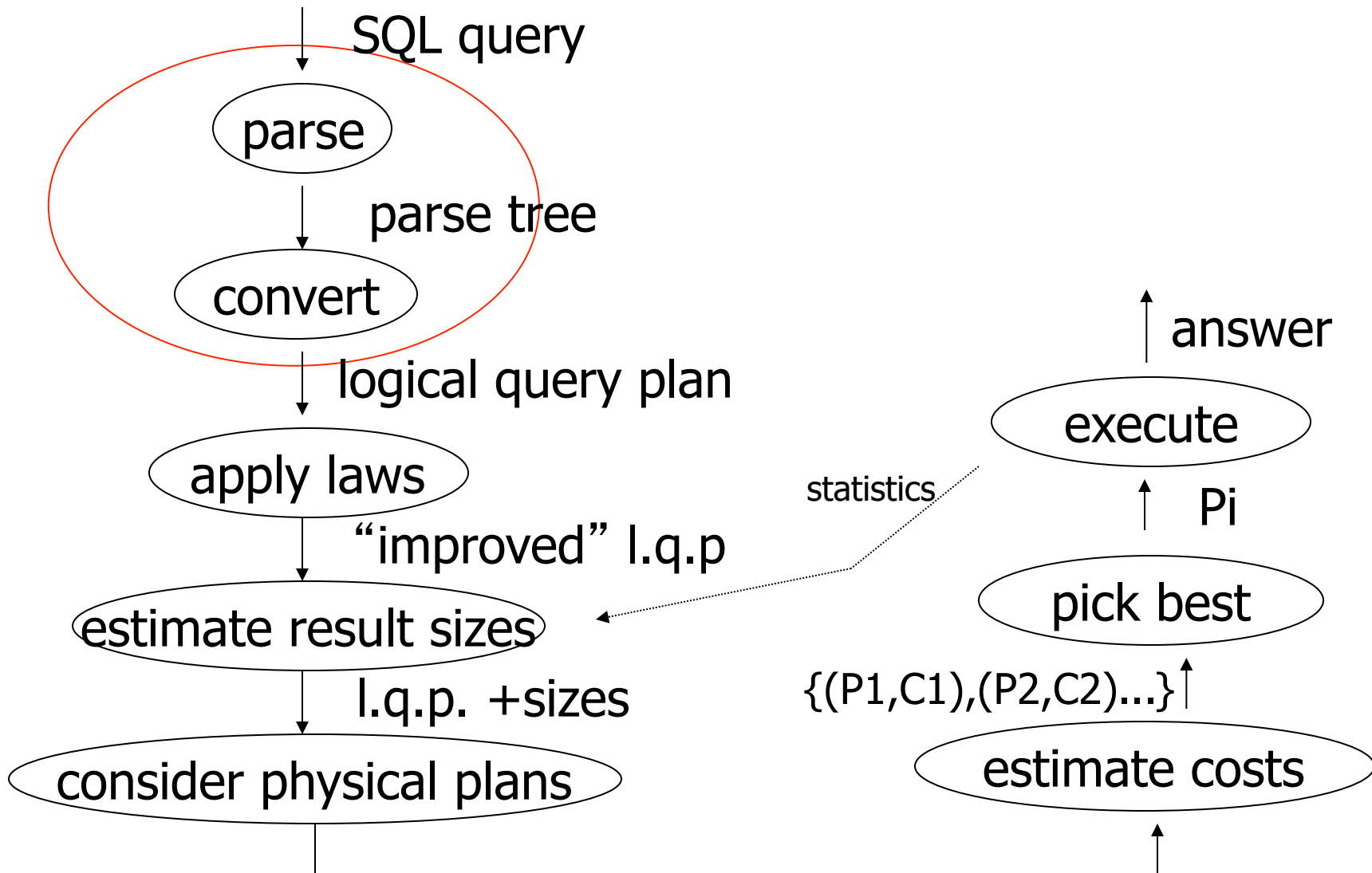
## 08: Query Processing Parsing and Analysis

Boris Glavic

Slides: adapted from a course taught by
Hector Garcia-Molina, Stanford InfoLab

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER
SCIENCE

SQL query

parse

parse tree

convert

logical query plan

apply laws

"improved" l.q.p

statistics

answer

execute

Pi

estimate result sizes

pick best

l.q.p. +sizes

{(P1,C1),(P2,C2)...}

consider physical plans

estimate costs

{P1,P2,.....}

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Parsing, Analysis, Conversion

## 1. Parsing
 – Transform SQL text into syntax tree

## 2. Analysis
 – Check for semantic correctness
 – Use database catalog
 – E.g., unfold views, lookup functions and attributes, check scopes

## 3. Conversion
 – Transform into internal representation
 – Relational algebra or QBM

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Analysis and Conversion

- Usually intertwined
- The internal representation is used to store analysis information
- Create an initial representation and complete during analysis

COMPUTER SCIENCE

Notes 8 - Parsing and Analysis

IIT College of
Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

# Parsing, Analysis, Conversion

1. Parsing
2. Analysis
3. Conversion

Notes 8 - Parsing and Analysis

IIT College of
Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Parsing

- SQL -> Parse Tree
- Covered in compiler courses and books
- Here only short overview

COMPUTER SCIENCE

IIT College of
Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

# SQL Standard

- Standardized language
  - 86, 89, 92, 99, 03, 06, 08, 11
- DBMS vendors developed their own dialects

COMPUTER SCIENCE

IIT College of Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY
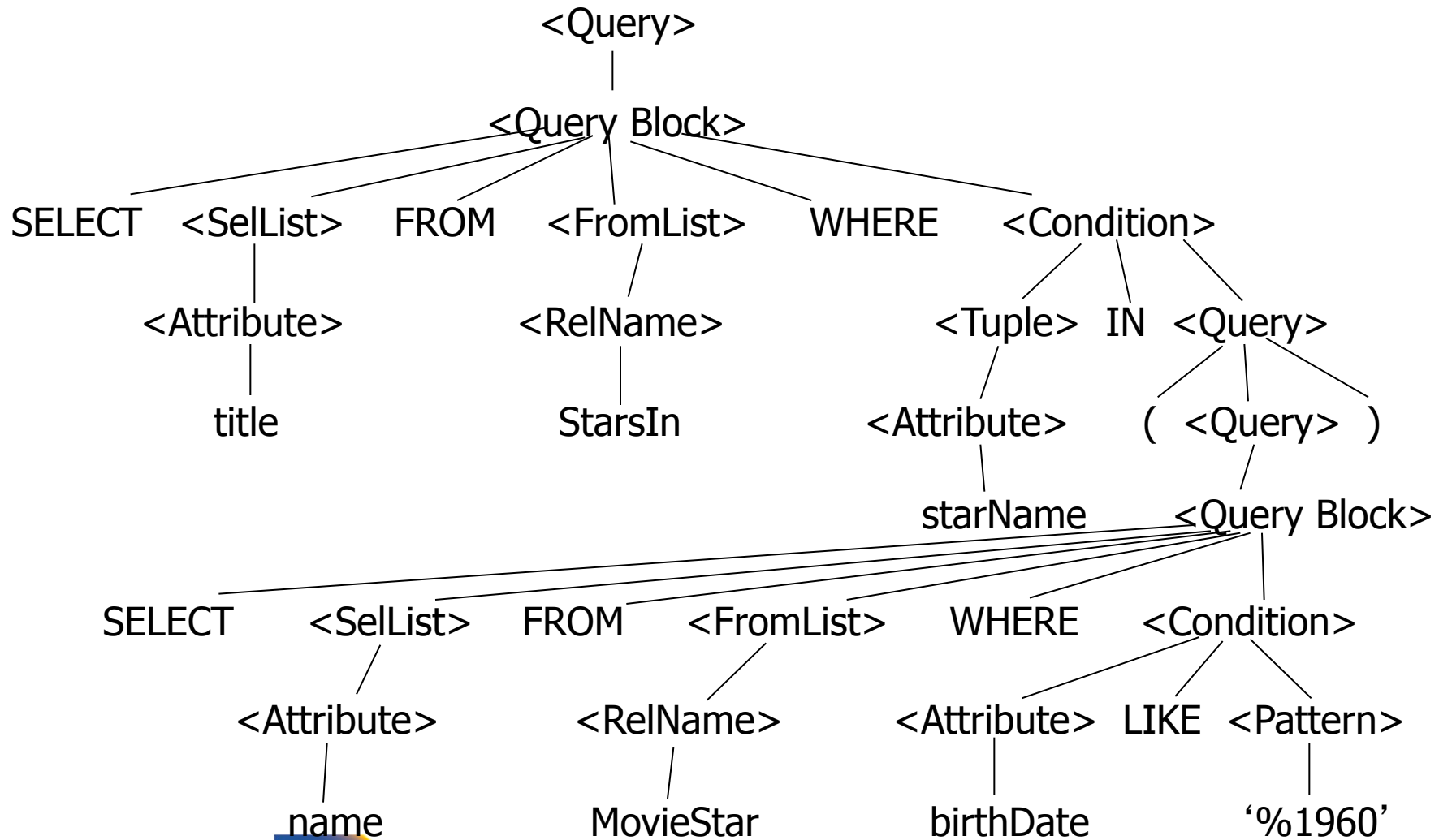
# Example:   SQL query

```
SELECT title
FROM StarsIn
WHERE starName IN (
        SELECT name
        FROM MovieStar
        WHERE birthdate LIKE '%1960'
);
```

(Find the movies with stars born in 1960)

Notes 8 - Parsing and Analysis     8

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER
SCIENCE

# Example: Parse Tree

<Query>
|
<Query Block>

SELECT   <SelList>   FROM   <FromList>   WHERE   <Condition>

<Attribute>          <RelName>                    <Tuple>  IN  <Query>

title                StarsIn              <Attribute>    (  <Query>  )

                                          starName      <Query Block>

SELECT   <SelList>   FROM   <FromList>   WHERE   <Condition>

<Attribute>          <RelName>            <Attribute>  LIKE  <Pattern>

name                 MovieStar            birthDate         '%1960'

# SQL Query Structure

- Organized in Query blocks

**SELECT** `<select_list>`

**FROM** `<from_list>`

**WHERE** `<where_condition>`

**GROUP BY** `<group_by_expressions>`

**HAVING** `<having_condition>`

**ORDER BY** `<order_by_expressions>`

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Query Blocks

- Only **SELECT** clause is mandatory
  - Some DBMS require **FROM**

**SELECT** (1 + 2) AS result

| result |
|--------|
| 3 |

COMPUTER SCIENCE

IIT College of
Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

# **SELECT** clause

- List of expressions and optional name assignment + optional **DISTINCT**
  - Attribute references: R.a, b
  - Constants: 1, 'hello', '2008-01-20'
  - Operators: (R.a + 3) * 2
  - Functions (maybe UDF): substr(R.a, 1,3)
    - Single result or **set functions**
  - Renaming: (R.a + 2) AS x

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# **SELECT** clause - example

```
SELECT substring(p.name,1,1) AS initial
          p.name
FROM person p
```

**person**

| name | gender |
|------|--------|
| Joe | male |
| Jim | male |

**result**

| initial | name |
|---------|------|
| J | Joe |
| J | Jim |

Notes 8 - Parsing and Analysis

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# **SELECT** clause – set functions

- Function `extrChar(string)`

**SELECT** `extrChar(p.name) AS n`
**FROM** `person p`

**person**

| name | gender |
|------|--------|
| Joe | male |
| Jim | male |

**result**

| n |
|---|
| J |
| o |
| e |
| J |
| i |
| m |

Notes 8 - Parsing and Analysis

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# SELECT clause – DISTINCT

```
SELECT DISTINCT gender
FROM person p
```

**person**

| name | gender |
|------|--------|
| Joe  | male   |
| Jim  | male   |

**result**

| gender |
|--------|
| male   |

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# FROM clause

- List of table expressions
  - Access to relations
  - Subqueries (need alias)
  - Join expressions
  - Table functions
  - Renaming of relations and columns

Notes 8 - Parsing and Analysis

IIT College of
Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# **FROM** clause examples

FROM R

-access table R

FROM R, S

-access tables R and S

FROM R JOIN S ON (R.a = S.b)

-join tables R and S on condition (R.a = S.b)

FROM R x

FROM R AS x

-Access table R and assign alias 'x'

COMPUTER
SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# **FROM** clause examples

```
FROM R x(c,d)
FROM R AS x(c,d)
        -using aliases x for R and c,d for its attribues
FROM (R JOIN S t ON (R.a = t.b)), T
        -join R and S, and access T
FROM (R JOIN S ON (R.a = S.b)) JOIN T
        -join tables R and S and result with T
FROM create_sequence(1,100) AS seq(a)
        -call table function
```

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# **FROM** clause examples

FROM

      `(SELECT count(*) FROM employee)`
      `AS empcnt(cnt)`

-count number of employee in subquery

Notes 8 - Parsing and Analysis

# **FROM** clause examples

```
SELECT *
FROM create_sequence(1,3) AS seq(a)
```

**result**

| a |
|---|
| 1 |
| 2 |
| 3 |

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# **FROM** clause examples

```
SELECT dep, headcnt
FROM (SELECT count(*) AS headcnt, dep
      FROM employee
      GROUP BY dep)
WHERE headcnt > 100
```

**employee**

| name | dep |
|------|-----|
| Joe | IT |
| Jim | Marketing |
| … | … |

**result**

| dep | headcnt |
|-----|---------|
| IT | 103 |
| Support | 2506 |
| … | … |

Notes 8 - Parsing and Analysis

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# **FROM** clause - correlation

- Correlation
  - Reference attributes from other FROM clause item
  - Attributes of $i^{th}$ entry only available in $j > i$
  - Semantics:
    - For each row in result of $i^{th}$ entry:
    - Substitute correlated attributes with value from current row and evaluate query

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Correlation - Example

```
SELECT name, chr
FROM employee AS e,
     extrChar(e.name) AS c(chr)
```

**result**

| name | chr |
|------|-----|
| Joe  | J   |
| Joe  | o   |
| Joe  | e   |
| Jim  | J   |
| Jim  | i   |
| ...  | ... |

**employee**

| name | dep       |
|------|-----------|
| Joe  | IT        |
| Jim  | Marketing |
| ...  | ...       |

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Correlation - Example

```
SELECT name
FROM (SELECT max(salary) maxsal
        FROM employee) AS m,
        (SELECT name
        FROM employee x
        WHERE x.salary = m.maxsal) AS e
```

**employee**

| name | salary |
|------|--------|
| Joe  | 20,000 |
| Jim  | 30,000 |
| ...  | ...    |

**result**

| name |
|------|
| Jim  |

# **WHERE** clause

- A condition
  - Attribute references
  - Constants
  - Operators (boolean)
  - Functions
  - Nested subquery expressions

- Result has to be boolean

COMPUTER SCIENCE

IIT College of
Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

# **WHERE** clause examples

```
WHERE R.a = 3
```
-comparison between attribute and constant
```
WHERE (R.a > 5) AND (R.a < 10)
```
-range query using boolean AND
```
WHERE R.a = S.b
```
-comparison between two attributes
```
WHERE (R.a * 2) > (S.b − 3)
```
-using operators

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Nested Subqueries

- Nesting a query within an expression
- Correlation allowed
  - Access FROM clause attributes
- Different types of nesting
  - Scalar subquery
  - Existential quantification
  - Universal quantification

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Nested Subqueries Semantics

- For each tuple produced by the FROM clause execute the subquery
  - If correlated attributes replace them with tuple values

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Scalar subquery

- Subquery that returns one result tuple
  - How to check?
  - -> Runtime error

```
SELECT *
FROM R
WHERE R.a = (SELECT count(*) FROM S)
```

Notes 8 - Parsing and Analysis

IIT College of
Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

# Existential Quantification

- `<expr> IN <subquery>`
  - Evaluates to true if <expr> equal to at least one of the results of the subquery

```
SELECT *
FROM users
WHERE name IN (SELECT name FROM
                        blacklist)
```

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Existential Quantification

- EXISTS `<subquery>`
  - Evaluates to true if <subquery> returns at least one tuple

```
SELECT *
FROM users u
WHERE EXISTS (SELECT * FROM
              blacklist b
          WHERE b.name = u.name)
```

IIT College of
Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Existential Quantification

- `<expr> <op> ANY <subquery>`
  - Evaluates to true if `<expr> <op> <tuple>` evaluates to true for **at least one** result tuple
  - Op is any comparison operator: =, <, >, …

```
SELECT *
FROM users
WHERE name = ANY (SELECT name FROM
                  blacklist)
```

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Universal Quantification

- `<expr> <op> ALL <subquery>`
  - Evaluates to true if `<expr> <op> <tuple>` evaluates to true for **all** result tuples
  - Op is any comparison operator: =, <, >, …

```
SELECT *
FROM nation
WHERE nname = ALL (SELECT nation FROM
                        blacklist)
```

Notes 8 - Parsing and Analysis

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Nested Subqueries Example

```
SELECT dep,name
FROM employee e
WHERE salary >= ALL (SELECT salary
                     FROM employee d
                     WHERE e.dep = d.dep)
```

**employee**

| name | dep | salary |
|---|---|---|
| Joe | IT | 2000 |
| Jim | IT | 300 |
| Bob | HR | 100 |
| Alice | HR | 10000 |
| Patrice | HR | 10000 |

**result**

| dep | Name |
|---|---|
| IT | Joe |
| HR | Alice |
| HR | Patrice |

COMPUTER SCIENCE

IIT College of
Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

# **GROUP BY** clause

- A list of expressions
  - Same as WHERE
  - No restriction to boolean
  - DBMS has to know how to compare = for data type

- Results are grouped by values of the expressions

- -> usually used for aggregation

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# **GROUP BY** restrictions

- If group-by is used then
  - SELECT clause can only use group by expressions or aggregation functions

**COMPUTER SCIENCE**

IIT College of
Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

# **GROUP BY** clause examples

```
GROUP BY R.a
```
    -group on single attribute
```
GROUP BY (1+2)
```
    -allowed but useless (single group)
```
GROUP BY salary / 1000
```
    -groups of salary values in buckets of 1000
```
GROUP BY R.a, R.b
```
    -group on two attributes

IIT College of
Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER
SCIENCE

```
SELECT count(*) AS numP,
      (SELECT count(*)
      FROM friends o
      WHERE o.with = f.name) AS numF
FROM (SELECT DISTINCT name FROM friends) f
GROUP BY (SELECT count(*)
            FROM friends o
            WHERE o.with = f.name)
```

**friends**

| name | with |
|------|------|
| Joe | Jim |
| Joe | Peter |
| Jim | Joe |
| Jim | Peter |
| Peter | Joe |

**result**

| numP | numF |
|------|------|
| 1 | 1 |
| 2 | 2 |

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# **HAVING** clause

- A boolean expression
- Applied after grouping and aggregation
  - Only references aggregation expressions and group by expressions

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# **HAVING** clause examples

…

HAVING `sum(R.a) > 100`

      -only return tuples with sum bigger than 100

…

`GROUP BY dep`

HAVING dep = 'IT' AND sum(salary) > 1000000

      -only return group 'IT' and sum threshold

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# **ORDER BY** clause

- A list of expressions
- Semantics: Order the result on these expressions

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# **ORDER BY** clause examples

```
ORDER BY R.a ASC
ORDER BY R.a
```
    -order ascending on R.a
```
ORDER BY R.a DESC
```
    -order descending on R.a
```
ORDER BY salary + bonus
```
    -order by sum of salary and bonus

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# New and Non-standard SQL features (excerpt)

- LIMIT / OFFSET
  - Only return a fix maximum number of rows
  - FETCH FIRST n ROWS ONLY (DB2)
  - row_number() (Oracle)

- Window functions
  - More flexible grouping
  - Return both aggregated results and input values

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Parsing, Analysis, Conversion

1. Parsing
2. Analysis
3. Conversion

Notes 8 - Parsing and Analysis

IIT College of
Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Analysis Goals

- Semantic checks
  - Table column exists
  - Operator, function exists
  - Determine type casts
  - Scope checks
- Rewriting
  - Unfolding views

Notes 8 - Parsing and Analysis

IIT College of
Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Semantic checks

```
SELECT *
FROM R
WHERE R.a + 3 > 5
```

- Table R exists?

- Expand ∗: which attributes in R?

- R.a is a column?

- Type of constants 3, 5?

- Operator + for types of R.a and 3 exists?

- Operator > for types of result of + and 5 exists?

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Database Catalog

- Stores information about database objects
- Aliases:
  - Information Schema
  - System tables
  - Data Dictionary

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Typical Catalog Information

- Tables
  - Name, attributes + data types, constraints
- Schema, DB
  - Hierarchical structuring of data
- Data types
  - Comparison operators
  - physical representation
  - Functions to (de)serialize to string

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Typical Catalog Information

- Functions (including aggregate/set)
  - Build-in
  - User defined (UDF)
- Triggers
- Stored Procedures
- ...

COMPUTER SCIENCE

IIT College of
Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

# Type Casts

- Similar to automatic type conversion in programming languages

- Expression: R.a + 3.0
  - Say R.a is of type integer
    - Search for a function +(int,float)
  - Does not exist?
    - Try to find a way to cast R.a, 3.0 or both to new data type
    - So that a function + exists for new types

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Scope checks

- Check that references are in correct scope

- E.g., if GROUP BY is present then SELECT clause expression can only reference group by expressions or aggregated values

COMPUTER SCIENCE

IIT College of
Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

# View Unfolding

- SQL allows for stored queries using CREATE VIEW

- Afterwards a view can be used in queries

- If view is not materialized, then need to replace view with its definition

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# View Unfolding Example

```
CREATE VIEW totalSalary AS
SELECT name, salary + bonus AS total
FROM employee


SELECT *
FROM totalSalary
WHERE total > 10000
```

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# View Unfolding Example

```
CREATE VIEW totalSalary AS
SELECT name, salary + bonus AS total
FROM employee


SELECT *
FROM (SELECT name,
             salary + bonus AS total
      FROM employee) AS totalSalary
WHERE total > 10000
```

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Analysis Summary

- Perform semantic checks
  - Catalog lookups (tables, functions, types)
  - Scope checks
- View unfolding
- Generate internal representation during analysis

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Parsing, Analysis, Conversion

1. Parsing
2. Analysis
3. Conversion

Notes 8 - Parsing and Analysis

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Conversion

- Create an internal representation
  - Should be useful for analysis
  - Should be useful optimization

- Internal representation
  - Relational algebra
  - Query tree/graph models
    - E.g., QGM (Query Graph Model) in Starburst

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Relational Alegbra

- Formal language
- Good for studying logical optimization and query equivalence (containment)
- Not informative enough for analysis
  - No datatype representation in algebra expressions
  - No meta-data

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Other Internal Representations

- Practical implementations
  - Mostly following structure of SQL query blocks
  - Store data type and meta-data (where necessary)

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Canonical Translation to Relational Algebra

- TEXTBOOK version of conversion
- Given an SQL query
- Return an equivalent relational algebra expression

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Relational Algebra Recap

- Formal query language
- Consists of operators
  - Input(s): relation
  - Output: relation
  - -> Composable
- Set and Bag semantics version

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

- Relation Schema
  - A set of attribute name-datatype pairs
- Relation (instance)
  - A (multi-)set of tuples with the same schema
- Tuple
  - List of attribute value pairs (or function from attribute name to value)

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER
SCIENCE

# Set- vs. Bag semantics

- Set semantics:
  - Relations are Sets
  - Used in most theoretical work

- Bag semantics
  - Relations are Multi-Sets
    - Each element (tuple) can appear more than once
  - SQL uses bag semantics

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER
SCIENCE

# Bag semantics notation

- We use $t^m$ to denote tuple t appears with multiplicity **m**

IIT College of
Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER
SCIENCE

# Set- vs. Bag semantics

Set

| Name | Purchase |
|------|----------|
| Peter | Guitar |
| Joe | Drum |
| Alice | Bass |

Bag

| Name | Purchase |
|------|----------|
| Peter | Guitar |
| Peter | Guitar |
| Joe | Drum |
| Alice | Bass |
| Alice | Bass |

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Operators

- Selection

- Renaming

- Projection

- Joins

  - Theta, natural, cross-product, outer, anti

- Aggregation

- Duplicate removal

- Set operations

# Selection

- Syntax: $\sigma_c (R)$
  - R is input
  - C is a condition

- Semantics:
  - Return all tuples that match condition C
  - Set: $\{ t \mid t \in R \text{ AND } t \text{ fulfills } C \}$
  - Bag: $\{ t^n \mid t^n \in R \text{ AND } t \text{ fulfills } C \}$

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Selection Example

- $\sigma_{a>5}$ (R)

R

| a | b |
|---|-----|
| 1 | 13 |
| 3 | 12 |
| 6 | 14 |

Result

| a | b |
|---|-----|
| 6 | 14 |

COMPUTER SCIENCE

Notes 8 - Parsing and Analysis

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Renaming

– Syntax: $\rho_A (R)$

- R is input
- A is list of attribute renamings b ← a

– Semantics:

- Applies renaming from A to inputs
- Set: $\{ t.A \mid t \, \varepsilon R \}$
- Bag: $\{ (t.A)^n \mid t^n \varepsilon R \}$

# Renaming Example

- $\rho_{c \leftarrow a} (R)$

R

| a | b |
|---|---|
| 1 | 13 |
| 3 | 12 |
| 6 | 14 |

Result

| c | b |
|---|---|
| 1 | 13 |
| 3 | 12 |
| 6 | 14 |

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Projection

– Syntax: $\Pi_A (R)$

- R is input
- A is list of projection expressions
- Standard: only attributes in A

– Semantics:

- Project all inputs on projection expressions
- Set: $\{ t.A \mid t \in R \}$
- Bag: $\{ (t.A)^n \mid t^n \in R \}$

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Projection Example

- $\Pi_b \ (R)$

R

| a | b |
|---|---|
| 1 | 13 |
| 3 | 12 |
| 6 | 14 |

Result

| b |
|---|
| 13 |
| 12 |
| 14 |

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Cross Product

- Syntax: R X S
  - R and S are inputs
- Semantics:
  - All combinations of tuples from R and S
  - = mathematical definition of cross product
  - Set: { (t,s) | t $\varepsilon$R AND s$\varepsilon$S }
  - Bag: { (t,s)$^{n*m}$ | t$^n\varepsilon$R AND s$^m\varepsilon$S }

Notes 8 - Parsing and Analysis

IIT College of
Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Cross Product Example

- R X S

R

| a | b |
|---|---|
| 1 | 13 |
| 3 | 12 |

S

| c | d |
|---|---|
| a | 5 |
| b | 3 |
| c | 4 |

Result

| a | b | c | d |
|---|---|---|---|
| 1 | 13 | a | 5 |
| 1 | 13 | b | 3 |
| 1 | 13 | c | 4 |
| 3 | 12 | a | 5 |
| 3 | 12 | b | 3 |
| 3 | 12 | c | 4 |

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Join

- Syntax: $R \bowtie_C S$
  - R and S are inputs
  - C is a condition

- Semantics:
  - All combinations of tuples from R and S that match C
  - Set: { (t,s) | t $\varepsilon$R AND s$\varepsilon$S AND (t,s) matches C}
  - Bag: { $(t,s)^{n*m}$ | $t^n\varepsilon$R AND $s^m\varepsilon$S AND (t,s) matches C}

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Join Example

- $R \bowtie_{a=d} S$

R

| a | b |
|---|----|
| 1 | 13 |
| 3 | 12 |

S

| c | d |
|---|---|
| a | 5 |
| b | 3 |
| c | 4 |

Result

| a | b | c | d |
|---|----|---|---|
| 3 | 12 | b | 3 |

Notes 8 - Parsing and Analysis

IIT College of
Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER
SCIENCE

# Natural Join

– Syntax: $R \bowtie S$

- R and S are inputs

– Semantics:

- All combinations of tuples from R and S that match on common attributes

- A = common attributes of R and S

- C = exclusive attributes of S

- Set: { (t,s.C) | t $\varepsilon$R AND s$\varepsilon$S AND t.A=s.A}

- Bag: { (t,s.C)$^{n*m}$ | t$^{n}\varepsilon$R AND s$^{m}\varepsilon$S AND t.A=s.A}

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Natural Join Example

- R ⋈ S

R

| a | b |
|---|---|
| 1 | 13 |
| 3 | 12 |

S

| c | a |
|---|---|
| a | 5 |
| b | 3 |
| c | 4 |

Result

| a | b | c |
|---|---|---|
| 3 | 12 | b |

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Left-outer Join

– Syntax: R $\bowtie_C$ S
  - R and S are inputs
  - C is condition

– Semantics:
  - R join S
  - t εR without match, fill S attributes with NULL

  { (t,s) | t εR AND sεS AND (t,s) matches C}

  union

  { (t, NULL(S)) | t εR AND NOT exists sεS: (t,s) matches C }

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Left-outer Join Example

- $R \bowtie_{a=d} S$

R

| a | b |
|---|---|
| 1 | 13 |
| 3 | 12 |

S

| c | d |
|---|---|
| a | 5 |
| b | 3 |
| c | 4 |

Result

| a | b | c | d |
|---|---|---|---|
| 1 | 13 | NULL | NULL |
| 3 | 12 | b | 3 |

Notes 8 - Parsing and Analysis

IIT College of Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Right-outer Join

– Syntax: R $\bowtie_C$ S

- R and S are inputs
- C is condition

– Semantics:

- R join S
- s εS without match, fill R attributes with NULL

{ (t,s) | t εR AND sεS AND (t,s) matches C}
union
{ (NULL(R),s) | s εS AND NOT exists tεR: (t,s) matches C }

Notes 8 - Parsing and Analysis

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Right-outer Join Example

- R $\bowtie_{a=d}$ S

R

| a | b |
|---|---|
| 1 | 13 |
| 3 | 12 |

S

| c | d |
|---|---|
| a | 5 |
| b | 3 |
| c | 4 |

Result

| a | b | c | d |
|------|------|---|---|
| NULL | NULL | a | 5 |
| 3 | 12 | b | 3 |
| NULL | NULL | c | 4 |

COMPUTER SCIENCE

IIT College of
Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

# Full-outer Join

– Syntax: R $\bowtie_C$ S

- R and S are inputs and C is condition

– Semantics:

{ (t,s) | t εR AND sεS AND (t,s) matches C}

union

{ (NULL(R),s) | s εS AND NOT exists tεR: (t,s) matches C }

union

{ (t, NULL(S)) | t εR AND NOT exists sεS: (t,s) matches C }

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER
SCIENCE

# Full-outer Join Example

- R ⋈$_{a=d}$ S

R

| a | b |
|---|----|
| 1 | 13 |
| 3 | 12 |

S

| c | d |
|---|---|
| a | 5 |
| b | 3 |
| c | 4 |

Result

| a | b | c | d |
|------|------|------|------|
| 1 | 13 | NULL | NULL |
| NULL | NULL | a | 5 |
| 3 | 12 | b | 3 |
| NULL | NULL | c | 4 |

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Semijoin

- Syntax: R ⋉ S  and R ⋊ S
    - R and S are inputs

- Semantics:
    - All tuples from R that have a matching tuple from relation S on the common attributes A

    { t | t εR AND exists sεS: t.A = s.A}

Notes 8 - Parsing and Analysis

# Semijoin Example

- R ⋉ S

R

| a | b |
|---|---|
| 1 | 13 |
| 3 | 12 |

S

| c | a |
|---|---|
| a | 5 |
| b | 3 |
| c | 4 |

Result

| a | b |
|---|---|
| 3 | 12 |

COMPUTER SCIENCE

Notes 8 - Parsing and Analysis

IIT College of Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

# Antijoin

- Syntax: R ▷ S
  - R and S are inputs

- Semantics:
  - All tuples from R that have no matching tuple from relation S on the common attributes A

  { t | t εR AND NOT exists sεS: t.A = s.A}

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Antijoin Example

- R ▷ S

R

| a | b |
|---|---|
| 1 | 13 |
| 3 | 12 |

S

| c | a |
|---|---|
| a | 5 |
| b | 3 |
| c | 4 |

Result

| a | b |
|---|---|
| 1 | 13 |

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Aggregation

– Syntax: $_G\alpha_A$ (R)

- A is list of aggregation functions
- G is list of group by attributes

– Semantics:

- Build groups of tuples according G and compute the aggregation functions from each group
- { (t.G, agg(G(t)) | t$\varepsilon$R }
- G(t) = { t' | t' $\varepsilon$R AND t'.G = t.G }

# Aggregation Example

- $_b\alpha_{sum(a)}$ (R)

R

| a | b |
|---|---|
| 1 | 1 |
| 3 | 1 |
| 6 | 2 |
| 3 | 2 |

Result

| sum(a) | b |
|--------|---|
| 4 | 1 |
| 9 | 2 |

COMPUTER SCIENCE

Notes 8 - Parsing and Analysis

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Duplicate Removal

- Syntax: $\delta(R)$
  - R is input
- Semantics:
  - Remove duplicates from input
  - Set: N/A
  - Bag: $\{ t^1 \mid t^n \varepsilon R \}$

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Duplicate Removal Example

- $\delta(R)$

R

| a | b |
|---|---|
| 1 | 13 |
| 1 | 13 |
| 6 | 14 |

Result

| a | b |
|---|---|
| 1 | 13 |
| 6 | 14 |

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Set operations

- Input: R and S
  - Have to have the same schema
    - Union compatible
  - Modulo attribute names
- Types
  - Union
  - Intersection
  - Set difference

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Union

- Syntax: R ∪ S
  - R and S are union-compatible inputs
- Semantics:
  - Set: { (t) | t εR OR tεS}
  - Bag: { $(t,s)^{n+m}$ | $t^n$εR AND $s^m$εS }
    - Assumption $t^n$ with $n < 1$ for tuple not in relation

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Union Example

- R ∪ S

R

| a |
|---|
| 1 |
| 3 |

S

| b |
|---|
| 1 |
| 2 |
| 3 |

Result

| a |
|---|
| 1 |
| 2 |
| 3 |
| 1 |
| 3 |

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Intersection

– Syntax: R ∩ S

  • R and S are union-compatible inputs

– Semantics:

  • Set: { (t) | t εR AND tεS}

  • Bag: { $(t,s)^{min(n,m)}$ | $t^n$εR AND $s^m$εS }

# Intersection Example

- R ∩ S

| R | | S | | Result |
|---|---|---|---|---|

R

| **a** |
|---|
| 1 |
| 3 |

S

| **b** |
|---|
| 1 |
| 2 |
| 3 |

Result

| **a** |
|---|
| 1 |
| 3 |

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Set Difference

- Syntax: R - S
  - R and S are union-compatible inputs

- Semantics:
  - Set: { (t) | t $\varepsilon$R AND NOT t$\varepsilon$S}
  - Bag: { $(t,s)^{n-m}$ | $t^n\varepsilon$R AND $s^m\varepsilon$S }

COMPUTER SCIENCE

IIT College of
Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

# Set Difference Example

- R - S

| R |
|---|
| **a** |
| 1 |
| 5 |

| S |
|---|
| **b** |
| 1 |
| 2 |
| 3 |

| Result |
|---|
| **a** |
| 5 |

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Canonical Translation to Relational Algebra

- TEXTBOOK version of conversion

- Given an SQL query

- Return an equivalent relational algebra expression

COMPUTER SCIENCE

IIT College of Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

# Canonical Translation

- **FROM** clause into joins and cross-products
  - Cross-product between list items
  - Joins into their algebra counter-part
- **WHERE** clause into selection
- **SELECT** clause into projection and renaming
  - If it has aggregation functions use aggreation
  - **DISTINCT** into duplicate removal

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Canonical Translation

- **GROUP BY** clause into aggregation
- **HAVING** clause into selection
- **ORDER BY** – no counter-part

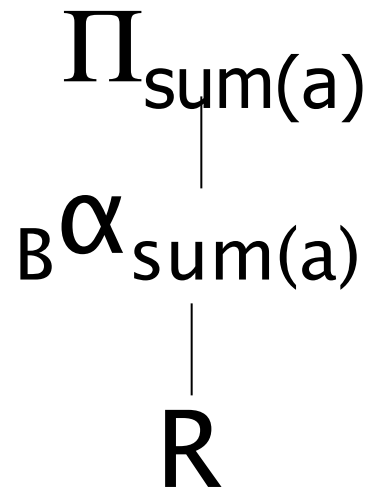- Then turn joins into crossproducts and selections

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

COMPUTER SCIENCE

# Set Operations

- **UNION ALL** into union
- **UNION** duplicate removal over union
- **INTERSECT ALL** into intersection
- **INTERSECT** add duplicate removal
- **EXCEPT ALL** into set difference
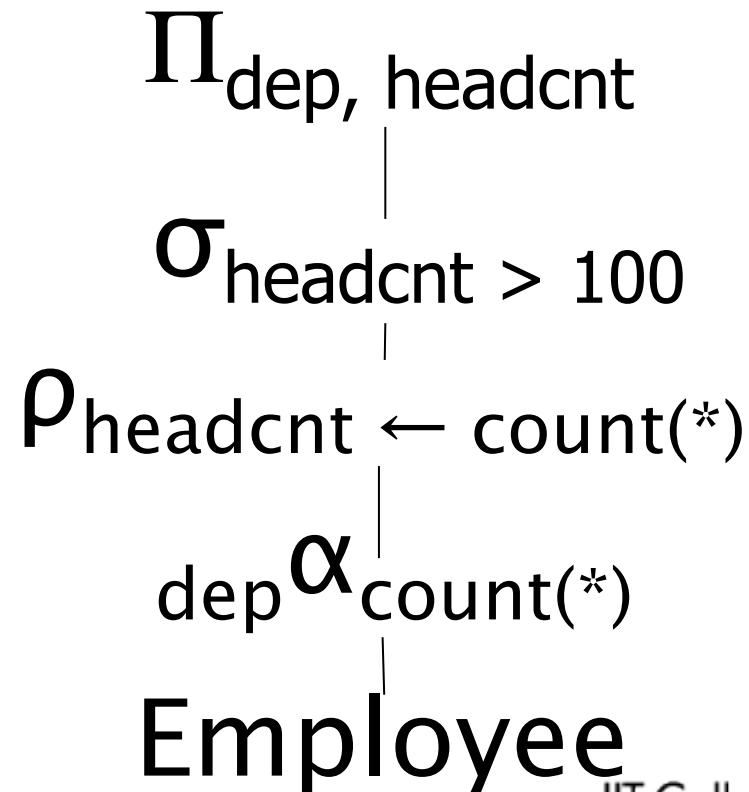- **EXCEPT** apply duplicate removal to inputs and then apply set difference

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Example:   Relational Algebra Translation

```
SELECT sum(R.a)
FROM R
GROUP BY b
```

$$\Pi_{\text{sum}(a)}$$

$$_B\alpha_{\text{sum}(a)}$$

$$R$$

COMPUTER SCIENCE

IIT College of
Science and Letters
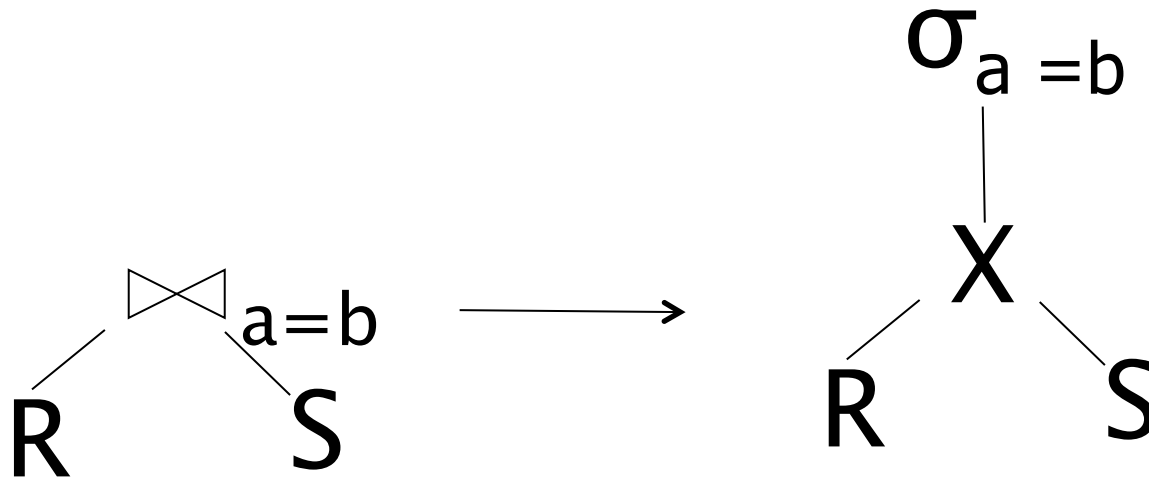
ILLINOIS INSTITUTE OF TECHNOLOGY

# Example: Relational Algebra Translation

```
SELECT dep, headcnt
FROM (SELECT count(*) AS headcnt, dep
        FROM employee
        GROUP BY dep)
WHERE headcnt > 100
```

$$\Pi_{dep, headcnt}$$

$$\sigma_{headcnt > 100}$$

$$\rho_{headcnt \leftarrow count(*)}$$

$$_{dep}\alpha_{count(*)}$$

$$Employee$$

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

# Example: Relational Algebra Translation

```
SELECT *
FROM R JOIN S ON (R.a = S.b)
```



$\bowtie_{a=b}$ over R, S $\longrightarrow$ $\sigma_{a=b}$ over X over R, S

COMPUTER SCIENCE

Notes 8 - Parsing and Analysis

IIT College of
Science and Letters

ILLINOIS INSTITUTE OF TECHNOLOGY

# Parsing and Analysis Summary

- SQL text -> Internal representation
- Semantic checks
- Database catalog
- View unfolding

COMPUTER SCIENCE

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY