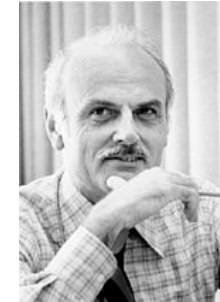# Relational Model

# Relational Model

- Many ad hoc models before 1970

  - Hard to work with

  - Hard to reason about

- **1970: Relational Model by Edgar Frank Codd**

  - Data are stored in relations (or tables)

  - Queried using a declarative language

  - DBMS converts declarative queries into procedural queries that are

    optimized and executed

- Key Advantages

  - Simple and clean mathematical model (based on logic)

  - Separation of declarative and procedural

# Relational Databases

| Flight | origin | destination | airline |
|--------|--------|-------------|---------|
| | VIE | LHR | BA |
| | LHR | EDI | BA |
| | LGW | GLA | U2 |
| | LCA | VIE | OS |

| Airport | code | city |
|---------|------|------|
| | VIE | Vienna |
| | LHR | London |
| | LGW | London |
| | LCA | Larnaca |
| | GLA | Glasgow |
| | EDI | Edinburgh |

Constants

VIE, LHR, …

BA, U2, …

Vienna, London, …

# Relational Databases

Relations

| Flight | origin | destination | airline |
|--------|--------|-------------|---------|
| VIE | LHR | BA |
| LHR | EDI | BA |
| LGW | GLA | U2 |
| LCA | VIE | OS |

| Airport | code | city |
|---------|------|------|
| VIE | Vienna |
| LHR | London |
| LGW | London |
| LCA | Larnaca |
| GLA | Glasgow |
| EDI | Edinburgh |

Constants

VIE, LHR, …

BA, U2, …

Vienna, London, …

# Relational Databases

**Relations**

| Flight | origin | destination | airline |
|--------|--------|-------------|---------|
|        | VIE    | LHR         | BA      |
|        | LHR    | EDI         | BA      |
|        | LGW    | GLA         | U2      |
|        | LCA    | VIE         | OS      |

Tuples

| Airport | code | city      |
|---------|------|-----------|
|         | VIE  | Vienna    |
|         | LHR  | London    |
|         | LGW  | London    |
|         | LCA  | Larnaca   |
|         | GLA  | Glasgow   |
|         | EDI  | Edinburgh |

**Constants**

VIE, LHR, …

BA, U2, …

Vienna, London, …

**Relational atoms**

Flight(LHR,EDI,BA)

Airport(LGW,London)

# Querying: Relational Algebra

**List all the airlines**

| Flight | origin | destination | airline |
|--------|--------|-------------|---------|
|  | VIE | LHR | BA |
|  | LHR | EDI | BA |
|  | LGW | GLA | U2 |
|  | LCA | VIE | OS |

| Airport | code | city |
|---------|------|------|
|  | VIE | Vienna |
|  | LHR | London |
|  | LGW | London |
|  | LCA | Larnaca |
|  | GLA | Glasgow |
|  | EDI | Edinburgh |

# Querying: Relational Algebra

**List all the airlines**

| Flight | origin | destination | airline |
|--------|--------|-------------|---------|
|        | VIE    | LHR         | BA      |
|        | LHR    | EDI         | BA      |
|        | LGW    | GLA         | U2      |
|        | LCA    | VIE         | OS      |

| Airport | code | city      |
|---------|------|-----------|
|         | VIE  | Vienna    |
|         | LHR  | London    |
|         | LGW  | London    |
|         | LCA  | Larnaca   |
|         | GLA  | Glasgow   |
|         | EDI  | Edinburgh |

{BA, U2, OS}

$\pi_{airline}$ Flight

# Querying: Relational Algebra

**List the codes of the airports in London**

| Flight | origin | destination | airline |
|---|---|---|---|
| | VIE | LHR | BA |
| | LHR | EDI | BA |
| | LGW | GLA | U2 |
| | LCA | VIE | OS |

| Airport | code | city |
|---|---|---|
| | VIE | Vienna |
| | LHR | London |
| | LGW | London |
| | LCA | Larnaca |
| | GLA | Glasgow |
| | EDI | Edinburgh |

# Querying: Relational Algebra

**List the codes of the airports in London**

| Flight | origin | destination | airline |
|--------|--------|-------------|---------|
| | VIE | LHR | BA |
| | LHR | EDI | BA |
| | LGW | GLA | U2 |
| | LCA | VIE | OS |

| Airport | code | city |
|---------|------|------|
| | VIE | Vienna |
| | LHR | London |
| | LGW | London |
| | LCA | Larnaca |
| | GLA | Glasgow |
| | EDI | Edinburgh |

{LHR, LGW}

$\pi_{code} \ (\sigma_{city='London'} \ Airport)$

# Querying: Relational Algebra

**List the airlines that fly directly from London to Glasgow**

| Flight | origin | destination | airline |
|--------|--------|-------------|---------|
|        | VIE    | LHR         | BA      |
|        | LHR    | EDI         | BA      |
|        | LGW    | GLA         | U2      |
|        | LCA    | VIE         | OS      |

| Airport | code | city      |
|---------|------|-----------|
|         | VIE  | Vienna    |
|         | LHR  | London    |
|         | LGW  | London    |
|         | LCA  | Larnaca   |
|         | GLA  | Glasgow   |
|         | EDI  | Edinburgh |

# Querying: Relational Algebra

**List the airlines that fly directly from London to Glasgow**

| Flight | origin | destination | airline |
|--------|--------|-------------|---------|
| | VIE | LHR | BA |
| | LHR | EDI | BA |
| | LGW | GLA | U2 |
| | LCA | VIE | OS |

| Airport | code | city |
|---------|------|------|
| | VIE | Vienna |
| | LHR | London |
| | LGW | London |
| | LCA | Larnaca |
| | GLA | Glasgow |
| | EDI | Edinburgh |

$\pi_{airline}$ ((Flight $\bowtie_{origin=code}$ ($\sigma_{city='London'}$ Airport)) $\bowtie_{destination=code}$ ($\sigma_{city='Glasgow'}$ Airport))

# Querying: Relational Algebra

**List the airlines that fly directly from London to Glasgow**

| Aux | origin | destination | airline | code | city | code | city |
|---|---|---|---|---|---|---|---|
| | LGW | GLA | U2 | LGW | London | GLA | Glasgow |

⬇

{U2}

$$\pi_{airline} ((Flight \bowtie_{origin=code} (\sigma_{city='London'} \; Airport)) \bowtie_{destination=code} (\sigma_{city='Glasgow'} \; Airport))$$

defines the auxiliary relation Aux

# Relational Algebra

- **Selection: σ**

- **Projection: π**

- **Cross product: ×**

- Natural join: ⋈

- **Rename: ρ**

- **Difference: \\**     in bold are the primitive operators

- **Union: ∪**

- Intersection: ∩

Formal definitions can be found in any database textbook

# Querying: Domain Relational Calculus

**List all the airlines**

| Flight | origin | destination | airline |
|--------|--------|-------------|---------|
|        | VIE    | LHR         | BA      |
|        | LHR    | EDI         | BA      |
|        | LGW    | GLA         | U2      |
|        | LCA    | VIE         | OS      |

| Airport | code | city |
|---------|------|------|
|         | VIE  | Vienna |
|         | LHR  | London |
|         | LGW  | London |
|         | LCA  | Larnaca |
|         | GLA  | Glasgow |
|         | EDI  | Edinburgh |

# Querying: Domain Relational Calculus

**List all the airlines**

| Flight | origin | destination | airline |
|--------|--------|-------------|---------|
|        | VIE    | LHR         | BA      |
|        | LHR    | EDI         | BA      |
|        | LGW    | GLA         | U2      |
|        | LCA    | VIE         | OS      |

| Airport | code | city      |
|---------|------|-----------|
|         | VIE  | Vienna    |
|         | LHR  | London    |
|         | LGW  | London    |
|         | LCA  | Larnaca   |
|         | GLA  | Glasgow   |
|         | EDI  | Edinburgh |

{BA, U2, OS}

{z | ∃x∃y Flight(x,y,z)}

# Querying: Domain Relational Calculus

**List the codes of the airports in London**

| Flight | origin | destination | airline |
|--------|--------|-------------|---------|
| | VIE | LHR | BA |
| | LHR | EDI | BA |
| | LGW | GLA | U2 |
| | LCA | VIE | OS |

| Airport | code | city |
|---------|------|------|
| | VIE | Vienna |
| | LHR | London |
| | LGW | London |
| | LCA | Larnaca |
| | GLA | Glasgow |
| | EDI | Edinburgh |

# Querying: Domain Relational Calculus

**List the codes of the airports in London**

| Flight | origin | destination | airline |
|--------|--------|-------------|---------|
|        | VIE    | LHR         | BA      |
|        | LHR    | EDI         | BA      |
|        | LGW    | GLA         | U2      |
|        | LCA    | VIE         | OS      |

| Airport | code | city      |
|---------|------|-----------|
|         | VIE  | Vienna    |
|         | LHR  | London    |
|         | LGW  | London    |
|         | LCA  | Larnaca   |
|         | GLA  | Glasgow   |
|         | EDI  | Edinburgh |

{LHR, LGW}

{x | ∃y Airport(x,y)  ∧  y = London}

# Querying: Domain Relational Calculus

**List the airlines that fly directly from London to Glasgow**

| Flight | origin | destination | airline |
|--------|--------|-------------|---------|
|        | VIE    | LHR         | BA      |
|        | LHR    | EDI         | BA      |
|        | LGW    | GLA         | U2      |
|        | LCA    | VIE         | OS      |

| Airport | code | city      |
|---------|------|-----------|
|         | VIE  | Vienna    |
|         | LHR  | London    |
|         | LGW  | London    |
|         | LCA  | Larnaca   |
|         | GLA  | Glasgow   |
|         | EDI  | Edinburgh |

# Querying: Domain Relational Calculus

**List the airlines that fly directly from London to Glasgow**

| Flight | origin | destination | airline |
|--------|--------|-------------|---------|
|        | VIE    | LHR         | BA      |
|        | LHR    | EDI         | BA      |
|        | LGW    | GLA         | U2      |
|        | LCA    | VIE         | OS      |

| Airport | code | city      |
|---------|------|-----------|
|         | VIE  | Vienna    |
|         | LHR  | London    |
|         | LGW  | London    |
|         | LCA  | Larnaca   |
|         | GLA  | Glasgow   |
|         | EDI  | Edinburgh |

⇩

{U2}

$\{z \mid \exists x \exists y \exists u \exists v \; Airport(x,u) \; \wedge \; u = London \; \wedge \; Airport(x,u) \; \wedge \; u = London\} \; \wedge \; Flight(x,y,z)\}$

# Domain Relational Calculus

$$\{x_1,...,x_k \mid \phi\}$$

first-order formula with

free variables $\{x_1,...,x_k\}$

But, we can express "**problematic**" queries, i.e., depend on the domain

$\{x \mid \forall y\, R(x,y)\}$     $\{x \mid \neg R(x)\}$     $\{x,y \mid R(x) \vee R(y)\}$

# Domain Relational Calculus

$$\{x_1,\ldots,x_k \mid \phi\}$$

first-order formula with

free variables $\{x_1,\ldots,x_k\}$

But, we can express "**problematic**" queries, i.e., depend on the domain

$\{x \mid \forall y\ R(x,y)\}$      $\{x \mid \neg R(x)\}$      $\{x,y \mid R(x) \vee R(y)\}$

**domain** = $\{1,2,3\}$

D = $\{R(1,1),\ R(1,2)\}$

**Ans** = { }

# Domain Relational Calculus

$$\{x_1, \ldots, x_k \mid \phi\}$$

first-order formula with

free variables $\{x_1, \ldots, x_k\}$

But, we can express "**problematic**" queries, i.e., depend on the domain

$\{x \mid \forall y\, R(x,y)\}$      $\{x \mid \neg R(x)\}$      $\{x,y \mid R(x) \lor R(y)\}$

**domain** = $\{1,2\}$

$D = \{R(1,1),\ R(1,2)\}$

**Ans** = $\{1\}$

# Domain Relational Calculus

$$\{x_1,...,x_k \mid \phi\}$$

first-order formula with

free variables $\{x_1,...,x_k\}$

But, we can express "**problematic**" queries, i.e., depend on the domain

$$\{x \mid \forall y\, R(x,y)\} \qquad \{x \mid \neg R(x)\} \qquad \{x,y \mid R(x) \lor R(y)\}$$

…thus, we adopt the active domain semantics - quantified variables range over

the active domain, i.e., the constants occurring in the input database

# Algebra = Calculus

A fundamental theorem (assuming the active domain semantics):

**Theorem:** The following query langauges are equally expressive

- Relational Algebra (**RA**)

- Domain Relational Calculus (**DRC**)

- Tuple Relational Calculus (**TRC**)

**Note:** Tuple relational calculus is the declarative language introduce by Codd. Domain relational calculus has been introduced later as a formalism closer to first-order logic

# Quiz!

**Is Glasgow reachable from Vienna?**

| Flight | origin | destination | airline |
|--------|--------|-------------|---------|
| | VIE | LHR | BA |
| | LHR | EDI | BA |
| | LGW | GLA | U2 |
| | LCA | VIE | OS |

| Airport | code | city |
|---------|------|------|
| | VIE | Vienna |
| | LHR | London |
| | LGW | London |
| | LCA | Larnaca |
| | GLA | Glasgow |
| | EDI | Edinburgh |



{ | ∃x∃y∃z∃w∃v  Airport(x,Vienna)  ∧  Airport(y,Glasgow)  ∧

Flight(x,z,w)  ∧  Flight(z,y,v)

**YES**

# Quiz!

**Is Glasgow reachable from Vienna?**

| Flight | origin | destination | airline |
|--------|--------|-------------|---------|
| | VIE | LHR | BA |
| | LHR | EDI | BA |
| | LGW | GLA | U2 |
| | LCA | VIE | OS |

| Airport | code | city |
|---------|------|------|
| | VIE | Vienna |
| | LHR | London |
| | LGW | London |
| | LCA | Larnaca |
| | GLA | Glasgow |
| | EDI | Edinburgh |



{ | ∃x∃y∃z∃w∃v Airport(x,Vienna) ∧ Airport(y,Glasgow) ∧

Flight(x,z,w) ∧ Flight(z,y,v)

**NO**

# Quiz!

**Is Glasgow reachable from Vienna?**

| Flight | origin | destination | airline |
|--------|--------|-------------|---------|
| | VIE | LHR | BA |
| | LHR | EDI | BA |
| | LGW | GLA | U2 |
| | LCA | VIE | OS |

| Airport | code | city |
|---------|------|------|
| | VIE | Vienna |
| | LHR | London |
| | LGW | London |
| | LCA | Larnaca |
| | GLA | Glasgow |
| | EDI | Edinburgh |



{ | ∃x∃y∃z∃w∃v Airport(x,Vienna) ∧ Airport(y,Glasgow) ∧

∃$z_1$∃$w_1$    Flight(x,z,w)    Flight(z,y,v)

∧ Flight(z,$z_1$,$w_1$) ∧              **YES**

# Quiz!

**Is Glasgow reachable from Vienna?**

| Flight | origin | destination | airline |
|--------|--------|-------------|---------|
| | VIE | LHR | BA |
| | LHR | EDI | BA |
| | LGW | GLA | U2 |
| | LCA | VIE | OS |

| Airport | code | city |
|---------|------|------|
| | VIE | Vienna |
| | LHR | London |
| | LGW | London |
| | LCA | Larnaca |
| | GLA | Glasgow |
| | EDI | Edinburgh |



$\{ \mid \exists x \exists y \exists z \exists w \exists v \; \text{Airport}(x, \text{Vienna}) \; \wedge \; \text{Airport}(y, \text{Glasgow}) \; \wedge$

$\exists z_1 \exists w_1 \quad \text{Flight}(x, z, w) \quad \text{Flight}(z, y, v)$

$\wedge \; \text{Flight}(z, z_1, w_1) \; \wedge$

**NO**

# Quiz!

**Is Glasgow reachable from Vienna?**

| Flight | origin | destination | airline |
|--------|--------|-------------|---------|
|        | VIE    | LHR         | BA      |
|        | LHR    | EDI         | BA      |
|        | LGW    | GLA         | U2      |
|        | LCA    | VIE         | OS      |

| Airport | code | city |
|---------|------|------|
|         | VIE  | Vienna |
|         | LHR  | London |
|         | LGW  | London |
|         | LCA  | Larnaca |
|         | GLA  | Glasgow |
|         | EDI  | Edinburgh |



Recursive query - not expressible in **RA/DRC/TRC**

(unless we bound the number of intermediate stops)

# Complexity of Query Languages

- The goal is to understand the complexity of evaluating a query over a database

- Our main technical tool is complexity theory

- What to measure? Queries may have a large output, and it would be unfair to count the output as "complexity"

- We therefore consider the following decision problems:
  - Query Output Tuple (QOT)
  - Boolean Query Evaluation (BQE)

# A Few Words on Complexity Theory

details can be found in the standard textbooks

see also notes on the webpage of the course

# Complexity Classes

Consider a function f : N → N

$$\text{TIME}(f(n)) \;=\; \{\Pi \mid \Pi \text{ is decided by some DTM in time } O(f(n))\}$$

$$\text{NTIME}(f(n)) \;=\; \{\Pi \mid \Pi \text{ is decided by some NTM in time } O(f(n))\}$$

$$\text{SPACE}(f(n)) \;=\; \{\Pi \mid \Pi \text{ is decided by some DTM using space } O(f(n))\}$$

$$\text{NSPACE}(f(n)) \;=\; \{\Pi \mid \Pi \text{ is decided by some NTM using space } O(f(n))\}$$

# Complexity Classes

- We can now recall the standard time and space complexity classes:

$$\text{PTIME} \quad = \quad \bigcup_{k>0} \text{TIME}(n^k)$$

$$\text{NP} \quad = \quad \bigcup_{k>0} \text{NTIME}(n^k)$$

$$\text{EXPTIME} \quad = \quad \bigcup_{k>0} \text{TIME}(2^{n^k})$$

$$\text{NEXPTIME} \quad = \quad \bigcup_{k>0} \text{NTIME}(2^{n^k})$$

$$\text{LOGSPACE} \quad = \quad \text{SPACE}(\log n)$$

$$\text{NLOGSPACE} \quad = \quad \text{NSPACE}(\log n)$$

these definitions are relying on two-tape Turing machines with a read-only and a read/write tape

$$\text{PSPACE} \quad = \quad \bigcup_{k>0} \text{SPACE}(n^k)$$

$$\text{EXPSPACE} \quad = \quad \bigcup_{k>0} \text{SPACE}(2^{n^k})$$

- For every complexity class C we can define its complementary class coC

# Relationship Among Complexity Classes

LOGSPACE $\subseteq$ NLOGSPACE $\subseteq$ PTIME $\subseteq$ NP, coNP $\subseteq$

PSPACE $\subseteq$ EXPTIME $\subseteq$ NEXPTIME, coNEXPTIME $\subseteq \cdots$

**Some useful notes:**

- For a deterministic complexity class C, coC = C

- coNLOGSPACE = NLOGSPACE

- It is generally believed that PTIME $\neq$ NP, but we don't know

- PTIME $\subset$ EXPTIME $\Rightarrow$ at least one containment between them is strict

- PSPACE = NPSPACE, EXPSPACE = NEXPSPACE, etc.

- But, we don't know whether LOGSPACE = NLOGSPACE

# Complete Problems

- These are the hardest problems in a complexity class

- A problem that is complete for a class C, it is unlikely to belong in a lower class

- A problem Π is complete for a complexity class C, or simply C-complete, if:

    1. Π ∈ C

    2. Π is C-hard, i.e., every problem Π' ∈ C can be efficiently reduced to Π

there exists a logspace algorithm that computes a function f such that

$$\mathbf{w} \in Π' \text{ iff } f(\mathbf{w}) \in Π \text{ - in this case we write } Π' \leq_L Π$$

- To show that Π is C-hard it suffices to reduce some C-hard problem Π' to it

# Some Complete Problems

- **NP-complete**

  - SAT (satisfiability of propositional formulas)

  - Many graph-theoretic problems (e.g., 3-colorability)

  - Traveling salesman

  - etc.

- **PSPACE-complete**

  - Quantified SAT (or simply QSAT)

  - Equivalence of two regular expressions

  - Many games (e.g., Geography)

  - etc.

# Back to Query Languages

# Complexity of Query Languages

- The goal is to understand the complexity of evaluating a query over a database

- Our main technical tool is complexity theory

- What to measure? Queries may have a large output, and it would be unfair to count the output as "complexity"

- We therefore consider the following decision problems:
  - Query Output Tuple (QOT)
  - Boolean Query Evaluation (BQE)

# Complexity of Query Languages

Some useful notation:

- Given a database D, and a query Q, Q(D) is the answer to Q over D

- **adom**(D) is the active domain of D - the constants occurring in D

- We write Q/k for the fact that the arity of Q is $k \geq 0$

**L** is some query language; for example, **RA**, **DRC**, etc. - we will see several query languages

---

QOT(**L**)

**Input:** a database D, a query $Q/k \in$ **L**, a tuple of constants $\mathbf{t} \in$ **adom**$(D)^k$

**Question:** $\mathbf{t} \in$ Q(D)?

# Complexity of Query Languages

Some useful notation:

- Given a database D, and a query Q, Q(D) is the answer to Q over D

- **adom**(D) is the active domain of D - the constants occurring in D

- We write Q/k for the fact that the arity of Q is k ≥ 0

**L** is some query language; for example, **RA**, **DRC**, etc. - we will see several query languages

---

BQE(**L**)

**Input:** a database D, a Boolean query Q ∈ **L**

**Question:** is Q(D) non-empty?

# Complexity of Query Languages

QOT(**L**)

**Input:** a database D, a query Q/k ∈ **L**, a tuple of constants **t** ∈ **adom**(D)$^k$

**Question:** **t** ∈ Q(D)?

BQE(**L**)

**Input:** a database D, a Boolean query Q ∈ **L**

**Question:** is Q(D) non-empty?

**Theorem:** QOT(**L**) $\equiv_L$ BQE(**L**), where **L** ∈ {**RA**, **DRC**, **TRC**}

($\equiv_L$ means logspace-equivalent)

# Complexity of Query Languages

(let us show this for domain relational calculus)

**Theorem:** QOT(**DRC**) $\equiv_L$ BQE(**DRC**)

**Proof:** ($\leq_L$) Consider a database $D$, a k-ary query $Q = \{x_1,\dots,x_k \mid \phi\}$, and a tuple $(t_1,\dots,t_k)$

Let $Q_{bool} = \{ \mid \exists x_1 \cdots \exists x_k (\phi \wedge x_1 = t_1 \wedge x_2 = t_2 \wedge \cdots \wedge x_k = t_k) \}$

Clearly, $(t_1,\dots,t_k) \in Q(D)$ iff $Q_{bool}(D)$ is non-empty

($\geq_L$) Trivial - a Boolean domain RC query is a domain RC query

# Complexity Measures

- **Combined complexity** - both D and Q are part of the input

- **Query complexity** - fixed D, input Q

BQE[D](**L**)

**Input:** a Boolean query Q ∈ **L**

**Question:** is Q(D) non-empty?

- **Data complexity** - input D, fixed Q

BQE[Q](**L**)

**Input:** a database D

**Question:** is Q(D) non-empty?

# Complexity of **RA**, **DRC**, **TRC**

**Theorem:** For **L** ∈ {**RA**, **DRC**, **TRC**} the following hold:

- BQE(**L**) is PSPACE-complete (combined complexity)

- BQE[D](**L**) is PSPACE-complete, for a fixed database D (query complexity)

- BQE[Q](**L**) is in LOGSPACE, for a fixed query Q ∈ **L** (data complexity)

**Proof hints:**

- Recursive algorithm that uses polynomial space in Q and logarithmic space in D

- Reduction from QSAT (a standard PSPACE-hard problem)

# Evaluating (Boolean) DRC Queries

Eval(D,φ)  -  for brevity we write φ instead of { | φ}

- **If** φ = $R(t_1,...,t_k)$, **then** YES iff $R(t_1,...,t_k) \in D$

- **If** φ = $\psi_1 \wedge \psi_2$, **then** YES iff Eval(D,$\psi_1$) = YES and Eval(D,$\psi_2$) = YES

- **If** φ = $\neg\psi$, **then** NO iff Eval(D,$\psi$) = YES

- **If** φ = $\exists x\ \psi(x)$, **then** YES iff for some t ∈ **adom**(D), Eval(D,$\psi(t)$) = YES

**Lemma:** It holds that

- Eval(D,φ) always terminates  -  this is trivial

- Eval(D,φ) = YES iff Q(D) is non-empty, where Q = { | φ}

- Eval(D,φ) uses $O(|φ| \cdot \log |φ| + |φ|^2 \cdot \log |D|)$ space

# Complexity of **RA**, **DRC**, **TRC**

**Theorem:** For **L** ∈ {**RA**, **DRC**, **TRC**} the following hold:

- BQE(**L**) is PSPACE-complete (combined complexity)

- BQE[D](**L**) is PSPACE-complete, for a fixed database D (query complexity)

- BQE[Q](**L**) is in LOGSPACE, for a fixed query Q ∈ **L** (data complexity)

**Proof hints:**

- Recursive algorithm that uses polynomial space in Q and logarithmic space in D

- Reduction from QSAT (a standard PSPACE-hard problem)

# Other Important Algorithmic Problems

SAT(**L**)

**Input:** a query Q ∈ **L**

**Question:** is there a (finite) database D such that Q(D) is non-empty?

---

EQUIV(**L**)

**Input:** two queries $Q_1$ ∈ **L** and $Q_2$ ∈ **L**

**Question:** $Q_1 \equiv Q_2$?  or  $Q_1(D) = Q_2(D)$ for every (finite) database D?

---

CONT(**L**)

**Input:** two queries $Q_1$ ∈ **L** and $Q_2$ ∈ **L**

**Question:** $Q_1 \subseteq Q_2$?  or  $Q_1(D) \subseteq Q_2(D)$ for every (finite) database D?

# Other Important Algorithmic Problems

SAT(**L**)

**Input:** a query Q ∈ **L**

**Question:** is there a (finite) database D such that Q(D) is non-empty?

EQUIV(**L**)

**Input:** two

**Question:**

**these problems are important**

**for optimization purposes**

abase D?

CONT(**L**)

**Input:** two queries $Q_1 \in$ **L** and $Q_2 \in$ **L**

**Question:** $Q_1 \subseteq Q_2$?  or  $Q_1(D) \subseteq Q_2(D)$ for every (finite) database D?

# Other Important Algorithmic Problems

SAT(**L**)

**Input:** a query Q ∈ **L**

**Question:** is there a (finite) database D such that Q(D) is non-empty?

- If the answer is no, then the input query Q makes no sense

- Query evaluation becomes trivial  -   the answer is always NO!

# Other Important Algorithmic Problems

EQUIV(**L**)

**Input:** two queries $Q_1 \in$ **L** and $Q_2 \in$ **L**

**Question:** $Q_1 \equiv Q_2$?  or  $Q_1(D) = Q_2(D)$ for every (finite) database D?

- Replace a query $Q_1$ with a query $Q_2$ that is easier to evaluate

- But, we have to be sure that $Q_1(D) = Q_2(D)$ for every database D

# Other Important Algorithmic Problems

CONT(**L**)

**Input:** two queries $Q_1 \in$ **L** and $Q_2 \in$ **L**

**Question:** $Q_1 \subseteq Q_2$?  or  $Q_1(D) \subseteq Q_2(D)$ for every (finite) database D?

- Approximate a query $Q_1$ with a query $Q_2$ that is easier to evaluate

- But, we have to be sure that $Q_1(D) \subseteq Q_2(D)$ for every database D

# SAT is Undecidable

**Theorem:** For **L** ∈ {**RA**, **DRC**, **TRC**}, SAT(**L**) is undecidable

**Proof hint:** By reduction from the halting problem.

Given a Turing machine M, we can construct a query $Q_M$ ∈ **L** such that:

M halts on the empty string   iff   there exists a database D such that Q(D) is non-empty



**Note:** Actually, this result goes back to the 1950 when

Boris A. Trakhtenbrot proved that the problem of deciding

whether a first-order sentence has a finite model is undecidable

# EQUIV and CONT are Undecidable

An easy consequence of the fact that SAT is undecidable is that:

**Theorem:** For **L** $\in$ {**RA**, **DRC**, **TRC**}, EQUIV(**L**) and CONT(**L**) are undecidable

**Proof:** By reduction from the complement of SAT(**L**)

- Consider a query Q $\in$ **L**  -  i.e., an instance of SAT(**L**)

- Let Q' be a query that is unsatisfiable, i.e., Q'(D) is empty for every D

- For example, when **L** = **DRC**, Q' can be the query { | $\exists$x R(x) $\wedge$ $\neg$R(x)}

- Clearly, Q is unsatisfiable  iff  Q $\equiv$ Q' (or even Q $\subseteq$ Q')

# Recap

- The main languages for querying relational databases are:

  - Relational Algebra (**RA**)

  - Domain Relational Calcuclus (**DRC**)

  - Tuple Relational Calculus (**TRC**)

$$\boxed{\textbf{RA} = \textbf{DRC} = \textbf{TRC}}$$

(under the active domain semantics)

- Evaluation is decidable, and highly tractable in data complexity

  - **Foundations of the database industry**

  - The core of SQL is equally expressive to **RA/DRC/TRC**

- Satisfiability, equivalence and containment are undecidable

  - **Perfect query optimization is impossible**