

Topic 1: Data models and query languages

Unit 2: Logic & relational calculus

Lecture 4

Wolfgang Gatterbauer

CS7240 Principles of scalable data management (sp23)

<https://northeastern-datalab.github.io/cs7240/sp23/>

1/20/2023

Topic 1: Data Models and Query Languages

- **Lecture 1 (Tue 1/10):** Course introduction / T1-U1 SQL / PostgreSQL setup / SQL Activities
- **Lecture 2 (Fri 1/13):** T1-U1 SQL
- **Lecture 3 (Tue 1/17):** T1-U1 SQL
- **Lecture 4 (Fri 1/20):** T1-U2 Logic & Relational Calculus
- **Lecture 5 (Tue 1/24):** T1-U1 Logic & Relational Calculus
- **Lecture 6 (Fri 1/27):** T1-U3 Relational Algebra & Codd's Theorem
- **Lecture 7 (Tue 1/31):** T1-U3 Relational Algebra & Codd's Theorem
- **Lecture 8 (Fri 2/3):** T1-U4 Datalog & Recursion
- **Lecture 9 (Tue 2/7):** T1-U4 Datalog & Recursion
- **Lecture 10 (Tue 2/10):** T1-U4 Datalog & Recursion

Pointers to relevant concepts & supplementary material:

- **Unit 1. SQL:** [SAMS'12], [CS 3200], [Cow'03] Ch3 & Ch5, [Complete'08] Ch6, [Silberschatz+'20] Ch3.8
- **Unit 2. Logic & Relational Calculus:** First-Order Logic (FOL), relational calculus (RC): [Barland+'08] 4.1.2 & 4.2.1 & 4.4, [Genesereth+] Ch6, [Halpern+'01], [Cow'03] Ch4.3 & 4.4, [Elmasri, Navathe'15] Ch8.6 & Ch8.7, [Silberschatz+'20] Ch27.1 & Ch27.2, [Alice'95] Ch3.1-3.3 & Ch4.2 & Ch4.4 & Ch5.3-5.4, [Barker-Plummer+'11] Ch11
- **Unit 3. Relational Algebra & Codd's Theorem:** Relational Algebra (RA), Codd's theorem: [Cow'03] Ch4.2, [Complete'08] Ch2.4 & Ch5.1-5.2, [Elmasri, Navathe'15] Ch8, [Silberschatz+'20] Ch2.6, [Alice'95] Ch4.4 & Ch5.4
- **Unit 4. Datalog & Recursion:** Datalog, recursion, Stratified Datalog with negation, Datalog evaluation strategies, Stable Model semantics, Answer Set Programming (ASP): [Complete'08] Ch5.3, [Cow'03] Ch 24, [Koutris'19] L9 & L10, [G., Suciu'10]
- **Unit 5. Alternative Data Models:** NoSQL: [Hellerstein, Stonebraker'05], [Sadalage, Fowler'12], [Harrison'16]

Queries and the connection to logic

- Why logic?
- A crash course in FOL
- Relational Calculus
 - Syntax and Semantics
 - Domain Independence and Safety

Logic as foundation of Computer Science and Databases

- Logic has had an immense impact on CS
- Computing has strongly driven a particular branch of logic: **finite model theory**
 - That is, **First-order logic (FOL)** restricted to finite models
 - Has strong connections to complexity theory
 - The basis of various branches in Artificial Intelligence (not the ones favored today)
- It is a natural tool to capture and attack fundamental problems in data management
 - Relations as first-class citizens
 - Inference for assuring data integrity (integrity constraints)
 - **Inference for question answering (queries)**
- It has been used for developing and analyzing the relational model from the early days [Codd'72]

Based on material by Benny Kimelfeld and Oded Shmueli for 236363 Database Management Systems, Technion, 2018.

See also: Halpern, Harper, Immerman, Kolaitis, Vardi, Vianu. "On the unusual effectiveness of logic in computer science", 2001. <https://doi.org/10.2307/2687775>

A play on: Wigner. "The unreasonable effectiveness of mathematics in the natural sciences", 1960. https://doi.org/10.1142/9789814503488_0018

Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>

Why has Logic turned out to be so powerful?

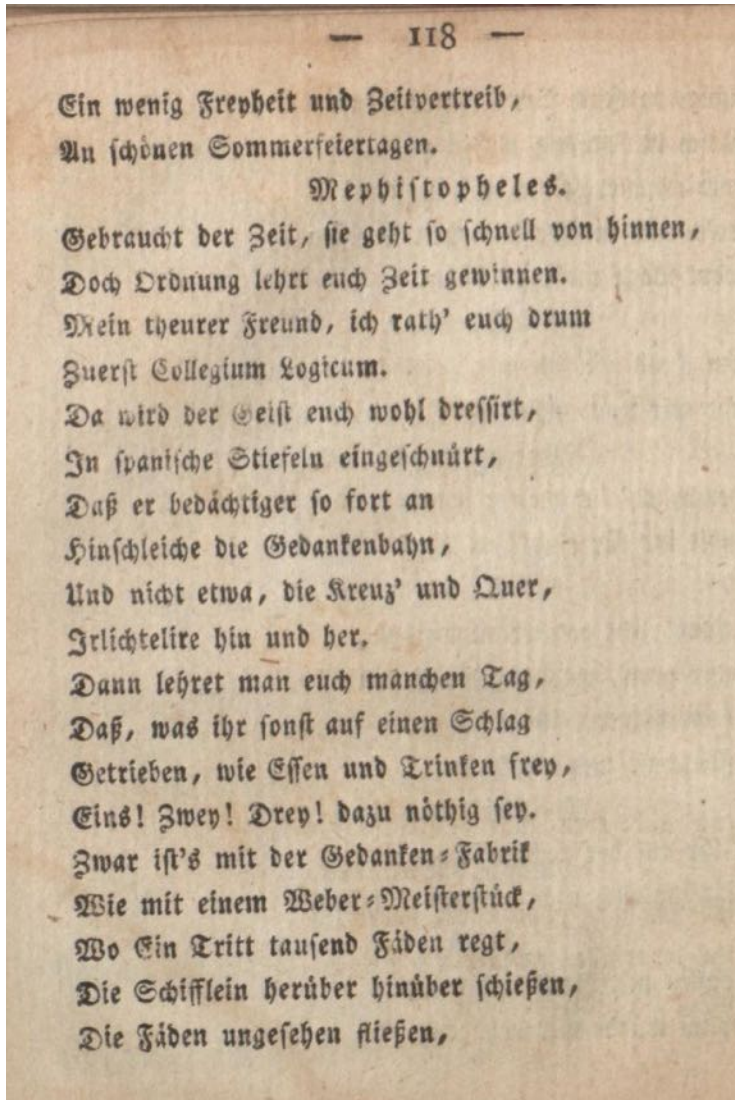
- Basic Question: What on earth does an obscure, old intellectual discipline have to do with the youngest intellectual discipline?
- Cosma R. Shalizi, CMU:
 - “If, in 1901, a talented and sympathetic outsider had been called upon (say, by a granting-giving agency) to survey the sciences and name the branch that would be least fruitful in century ahead, his choice might well have settled upon **mathematical logic**, an exceedingly recondite field whose practitioners could all have fit into a small auditorium. It had no practical applications, and not even that much mathematics to show for itself: its crown was an exceedingly obscure definition of cardinal numbers.”

See here for pointers to some of these discussions:

https://en.wikipedia.org/wiki/Cardinal_number



Logics as the start of everything ["Mephistopheles" 1806]



MEPHISTOPHELES

Gebraucht der Zeit, sie geht so schnell von hinnen,
Doch Ordnung lehrt Euch Zeit gewinnen.
Mein teurer Freund, ich rat Euch drum
Zuerst Collegium Logicum.

Da wird der Geist Euch wohl dressiert,
In spanische Stiefeln eingeschnürt,
Daß er bedächtiger so fortan
Hinschleiche die Gedankenbahn,
Und nicht etwa, die Kreuz und Quer,
Irrlichteliere hin und her.

...

MEPHISTOPHELES

Use your time well: it slips away so fast, yet
Discipline will teach you how to win it.
My dear friend, I'd advise, in sum,
First, the Collegium Logicum.

There your mind will be trained,
As if in Spanish boots, constrained,
So that painfully, as it ought,
It creeps along the way of thought,
Not flitting about all over,
Wandering here and there.

...

Source: Johan Wolfgang von Goethe. Faust Part I: Scene IV: The Study. ~1806. https://www.deutschestextarchiv.de/book/view/goethe_faust01_1808?p=124,

English Translation: <https://www.poetryintranslation.com/PITBR/German/FaustI/ScenesIVtoVI.php>

Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>

Back to The Future

- M. Davis (1988): Influences of Mathematical Logic on Computer Science:
 - “When I was a student, even the topologists regarded mathematical logicians as living in outer space. Today the connections between logic and computers are a matter of engineering practice at every level of computer organization.”
- Question: Why on earth?

Birth of Computer Science: 1930s

- Church, Gödel, Kleene, Post, Turing: Mathematical proofs have to be “machine checkable” - computation lies at the heart of mathematics!
 - Fundamental Question: What is “machine checkable”?
- Fundamental Concepts:
 - algorithm: a procedure for solving a problem by carrying out a precisely determined sequence of simpler, unambiguous steps
 - distinction between hardware and software
 - a universal machine: a machine that can execute arbitrary programs
 - a programming language: notation to describe algorithms

Leibniz's Dream

An Amazing Dream: a universal mathematical language, *lingua characteristica universalis*, in which all human knowledge can be expressed, and calculational rules, *calculus ratiocinator*, carried out by machines, to derive all logical relationships

- “If controversies were to arise, there would be no more need of disputation between two philosophers than between two accountants. For it would suffice to take their pencils in their hands, and say to each other: **Calculemus**—Let us calculate.”

Example: Aristotle' Syllogisms



- “All humans are mortal”



Example: Aristotle' Syllogisms



- “All humans are mortal”
- “For all x , if x is a human, then x is mortal”



Example: Aristotle' Syllogisms



- “All humans are mortal”
- “For all x, if x is a human, then x is mortal”

- $\forall x [\text{Human}(x) \rightarrow \text{Mortal}(x)]$ *Do you see the connection to referential integrity constraints ?*

Product

<u>PName</u>	Price	Category	cid
Gizmo	\$19.99	Gadgets	1
Powergizmo	\$29.99	Gadgets	1
SingleTouch	\$14.99	Photography	2
MultiTouch	\$203.99	Household	3

Company

<u>cid</u>	CName	StockPrice	Country
1	GizmoWorks	25	USA
2	Canon	65	Japan
3	Hitachi	15	Japan

Example: Aristotle' Syllogisms

- “All humans are mortal”
- “For all x, if x is a human, then x is mortal”
- $\forall x [\text{Human}(x) \rightarrow \text{Mortal}(x)]$ *Do you see the connection to referential integrity constraints*

$\forall x [\text{Product}(_,_,_,x) \rightarrow \text{Company}(x,_,_,_)]$

Product

<u>PName</u>	Price	Category	cid
Gizmo	\$19.99	Gadgets	1
Powergizmo	\$29.99	Gadgets	1
SingleTouch	\$14.99	Photography	2
MultiTouch	\$203.99	Household	3

Company

<u>cid</u>	CName	StockPrice	Country
1	GizmoWorks	25	USA
2	Canon	65	Japan
3	Hitachi	15	Japan

Logic and Databases

Two main uses of logic in databases:

- Logic used as a **database query language** to express questions asked against databases (our main focus)
- Logic used as a specification language to express **integrity constraints** in databases (product/company example from previous slide)

Why Logic?

- Logic provides both a unifying framework and a set of tools for formalizing and studying data management tasks.

Logic in Computer Science

- During the past fifty years there has been extensive, continuous, and growing interaction between logic and computer science. In many respects, logic provides computer science with both **a unifying foundational framework** and a tool for modeling computational systems. In fact, logic has been called “the calculus of computer science”.
- The argument is that logic plays a fundamental role in computer science, similar to that played by calculus in the physical sciences and traditional engineering disciplines.
 - Indeed, logic plays an important role in areas of computer science as disparate as machine architecture, computer-aided design, programming languages, databases, artificial intelligence, algorithms, and computability and complexity.

Queries and the connection to logic

- Why logic?
- A crash course in FOL
- Relational Calculus
 - Syntax and Semantics
 - Domain Independence and Safety

First-Order Logic: some notions

- Objects, e.g., 2
- Predicates (relationships), e.g., $2 < 3$
 - notice predicates are Boolean functions
 - e.g., Define $f(x,y)=\text{true}$ iff $x < y$. Thus $f(2,3)=\text{true}$
- Operations (non-Boolean functions), e.g., $2 + 3$
 - such functions usually return an object from the same domain as the inputs
- Logical operations, e.g., “and” (\wedge), “or” (\vee), “implies” (\rightarrow)
 - Both inputs and outputs are Boolean
- Quantifiers, e.g., “for all” (\forall), “exists” (\exists)

First-Order Logic

- A formalism for specifying properties of mathematical structures, such as graphs, partial orders, groups, rings, fields, . . .
- For any given structure, we can verify whether the properties hold

- **Mathematical Structure:**

- $A = (D, R_1, \dots, R_k, f_1, \dots, f_l)$

- D is a non-empty set: universe, or domain

$$D = \{1,2,3,4\}$$

R1

A	B	C
1	1	1
2	4	3

- R_i is an m -ary relation on D , for some m (i.e., $R_i \subseteq D^m$)

$$D^m \rightarrow \{T, F\}$$

$$\subseteq D \times D \times D$$

- f_j is an n -ary function on D , for some n (i.e., $f_i: D^n \rightarrow D$)

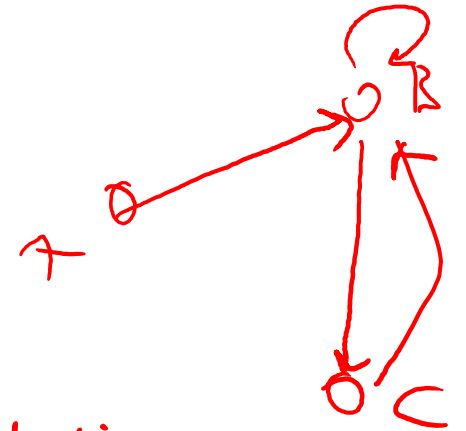
$$4 \cdot 4 \cdot 4 = 64$$

$$f(w_1, w_2) = w_1 + w_2$$

First-Order Logic on Graphs

Syntax:

- First-order variables: x, y, z, \dots (range over nodes)
- Atomic formulas: $E(x, y), x = y$
- Formulas:
 - Atomic Formulas, and
 - Boolean Connectives (\vee, \wedge, \neg), and
 - First-Order Quantifiers ($\exists x, \forall x$)



binary relation

Edge

from	to
A	B
B	B
C	B
B	C

Example properties for graph



- “node x has at least two distinct neighbors”



Assume schema is $E(\text{source}, \text{target})$, yet undirected. Thus for every edge $E(x,y)$, we also have $E(y,x)$. As convention later, we use $A(x,y)$ for directed edges = arcs

- “each node has at least two distinct neighbors”



Example properties for graph



- “node x has at least two distinct neighbors”

- $\exists y \exists z [E(\overset{\text{Alice}}{\cancel{x}}, y) \wedge E(\overset{\text{Alice}}{\cancel{x}}, z) \wedge y \neq z]$

- Notice: x is **free** in the above formula, which expresses a property of a node ~~x~~ .

- You can also think about this as a **query** (find nodes x that have ...)

Assume schema is $E(\text{source}, \text{target})$, yet undirected. Thus for every edge $E(x,y)$, we also have $E(y,x)$. As convention later, we use $A(x,y)$ for directed edges = arcs

Alice

- “each node has at least two distinct neighbors”



Example properties for graph



- “node x has at least two distinct neighbors”

- $\exists y \exists z [E(x, y) \wedge E(x, z) \wedge y \neq z]$

- Notice: x is **free** in the above formula, which expresses a property of a node ' x '.

- You can also think about this as a **query** (find nodes x that have ...)

Assume schema is $E(\text{source}, \text{target})$, yet undirected. Thus for every edge $E(x, y)$, we also have $E(y, x)$. As convention later, we use $A(x, y)$ for directed edges = arcs

- “each node has at least two distinct neighbors”

- $\forall x \exists y \exists z [E(x, y) \wedge E(x, z) \wedge y \neq z]$

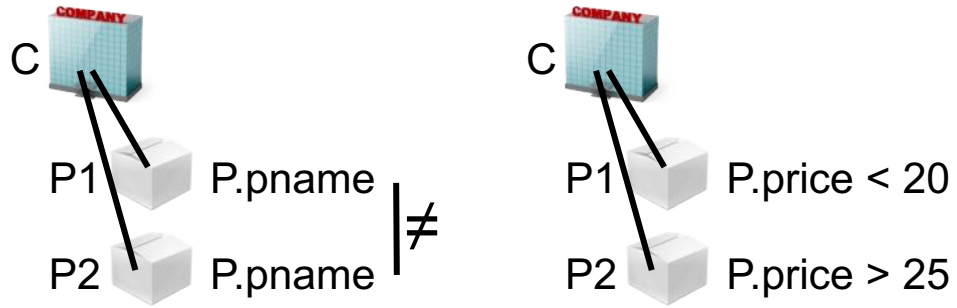
- The above is a **sentence**, that is, a formula with **no free variables**; it expresses a property of graphs.

We will sometimes use $\exists x, y, z$ as short form for $\exists x \exists y \exists z$



- “Find nodes that have at least two distinct neighbors” (query)

– $\{x \mid \exists y \exists z [E(x, y) \wedge E(x, z) \wedge y \neq z]\}$



- “each node has at least two distinct neighbors” (statement = Boolean query)

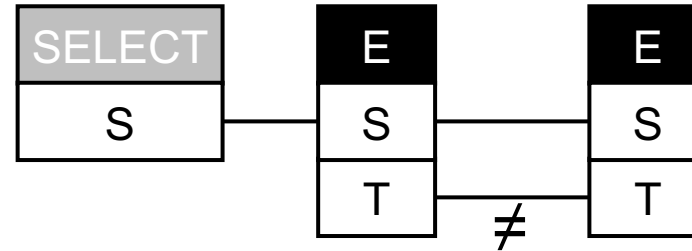
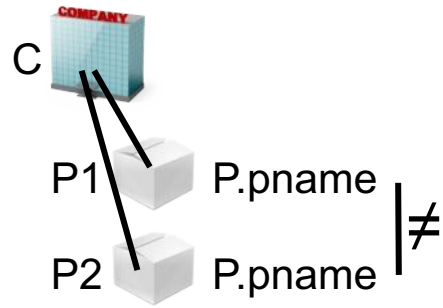
– $\forall x \exists y \exists z [E(x, y) \wedge E(x, z) \wedge y \neq z]$



- “Find nodes that have at least two distinct neighbors” (query)

– $\{x \mid \exists y \exists z [E(x, y) \wedge E(x, z) \wedge y \neq z]\}$

```
SELECT DISTINCT E1.S
FROM   E E1, E E2
WHERE  E1.S = E2.S
AND    E1.T <> E2.T
```



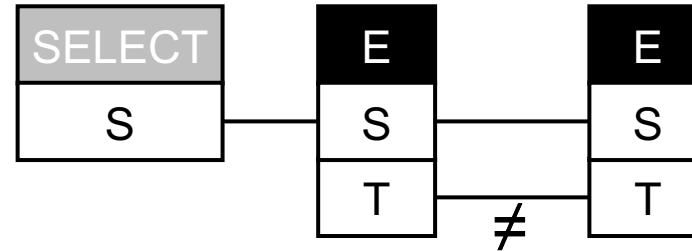
- “each node has at least two distinct neighbors” (statement = Boolean query)

– $\forall x \exists y \exists z [E(x, y) \wedge E(x, z) \wedge y \neq z]$
– $\neg(\exists x \neg(\exists y \exists z [E(x, y) \wedge E(x, z) \wedge y \neq z]))$



- “Find nodes that have at least two distinct neighbors” (query)

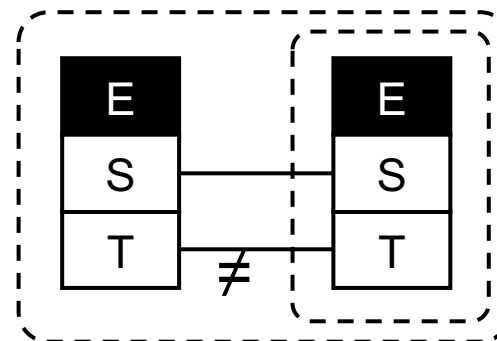
– $\{x \mid \exists y \exists z [E(x, y) \wedge E(x, z) \wedge y \neq z]\}$



```
SELECT DISTINCT E1.S
FROM   E E1, E E2
WHERE  E1.S = E2.S
AND    E1.T <> E2.T
```

- “each node has at least two distinct neighbors”

– $\forall x \exists y \exists z [E(x, y) \wedge E(x, z) \wedge y \neq z]$
– $\neg(\exists x \neg(\exists y \exists z [E(x, y) \wedge E(x, z) \wedge y \neq z]))$

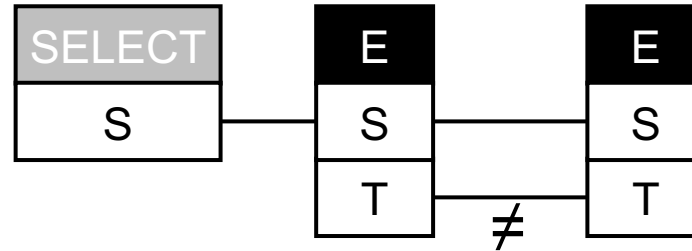
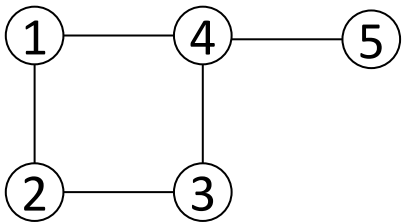


```
SELECT not exists
(SELECT *
FROM E E1
WHERE not exists
(SELECT *
FROM E E2
WHERE E1.S = E2.S
AND E1.T <> E2.T))
```



- “Find nodes that have at least two distinct neighbors” (query)

– $\{x \mid \exists y \exists z [E(x, y) \wedge E(x, z) \wedge y \neq z]\}$

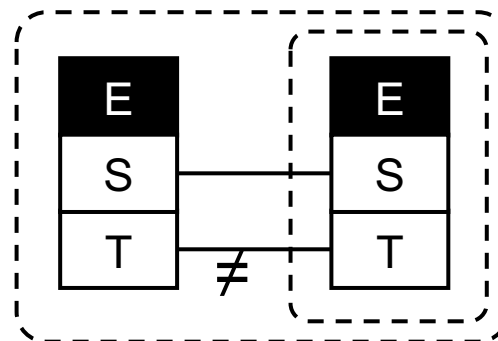
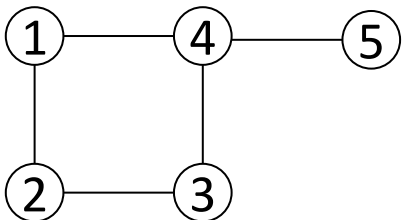


```
SELECT DISTINCT E1.S
FROM   E E1, E E2
WHERE  E1.S = E2.S
AND    E1.T <> E2.T
```



- “each node has at least two distinct neighbors”

– $\forall x \exists y \exists z [E(x, y) \wedge E(x, z) \wedge y \neq z]$
– $\neg(\exists x \neg(\exists y \exists z [E(x, y) \wedge E(x, z) \wedge y \neq z]))$



```
SELECT not exists
( SELECT *
FROM E E1
WHERE not exists
( SELECT *
FROM E E2
WHERE E1.S = E2.S
AND E1.T <> E2.T ))
```



Now in SQL

What is a minimal change to the two queries to evaluate them only over nodes 1-4



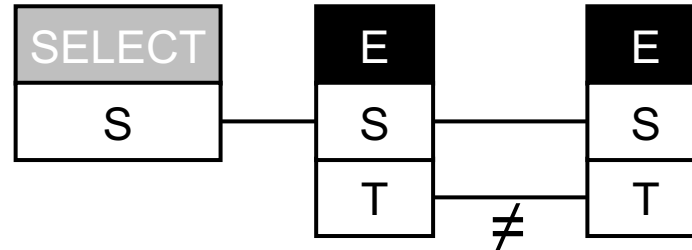
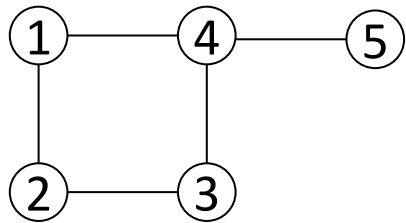
E(S,T)



501

- “Find nodes that have at least two distinct neighbors” (query)

– $\{x \mid \exists y \exists z [E(x, y) \wedge E(x, z) \wedge y \neq z]\}$

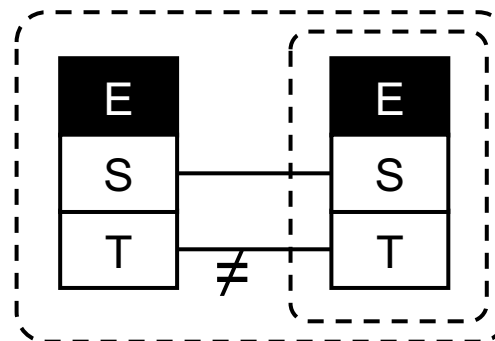
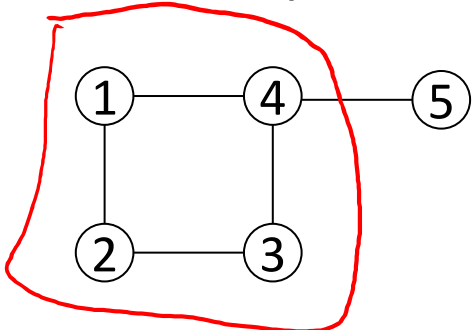


```
SELECT DISTINCT E1.S
FROM E E1, E E2
WHERE E1.S = E2.S
AND E1.T <> E2.T
```

{1, 2, 3, 4}

- “each node has at least two distinct neighbors”

– $\forall x \exists y \exists z [E(x, y) \wedge E(x, z) \wedge y \neq z]$
 – $\neg(\exists x \neg(\exists y \exists z [E(x, y) \wedge E(x, z) \wedge y \neq z]))$



```
SELECT not exists
(SELECT *
FROM E E1
WHERE not exists
(SELECT *
FROM E E2
WHERE E1.S = E2.S
AND E1.T <> E2.T))
```

FALSE

Now in SQL

A minimal change to the two queries to evaluate them only over nodes 1-4:

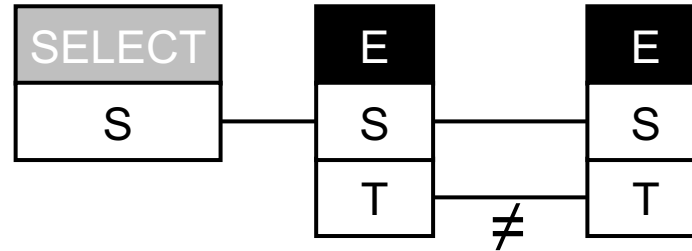
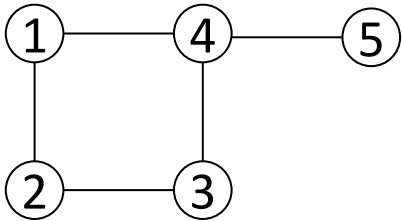
E(S,T)



501

- “Find nodes that have at least two distinct neighbors” (query)

– $\{x \mid \exists y \exists z [E(x, y) \wedge E(x, z) \wedge y \neq z]\}$

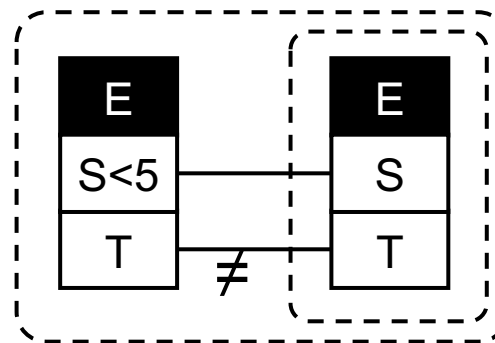
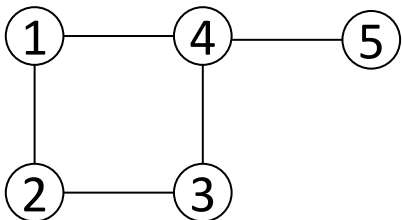


```
SELECT DISTINCT E1.S
FROM   E E1, E E2
WHERE  E1.S = E2.S
AND    E1.T <> E2.T
AND    E1.S < 5
```

{1, 2, 3, 4}

- “each node has at least two distinct neighbors”

– $\forall x \exists y \exists z [E(x, y) \wedge E(x, z) \wedge y \neq z]$
– $\neg(\exists x \neg(\exists y \exists z [E(x, y) \wedge E(x, z) \wedge y \neq z]))$



```
SELECT not exists
  (SELECT *
   FROM E E1
   WHERE E1.S < 5
   AND not exists
     (SELECT *
      FROM E E2
      WHERE E1.S = E2.S
      AND E1.T <> E2.T))
```

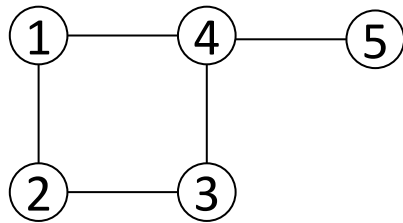
TRUE

Now in SQL with grouping

E(S,T)



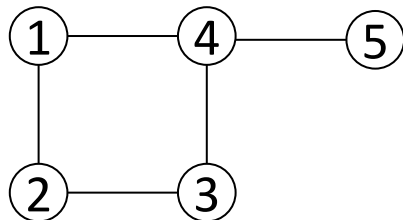
- “Find nodes that have at least two distinct neighbors” (query)



```
SELECT DISTINCT E1.S
FROM   E E1, E E2
WHERE  E1.S = E2.S
AND    E1.T <> E2.T
```

{1, 2, 3, 4}

- “each node has at least two distinct neighbors”



```
SELECT not exists
  (SELECT *
   FROM E E1
   WHERE not exists
     (SELECT *
      FROM E E2
      WHERE E1.S = E2.S
            AND E1.T <> E2.T))
```

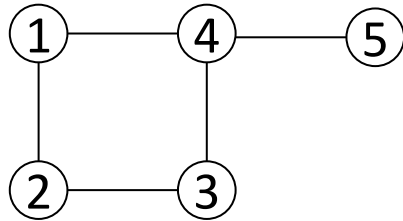
FALSE

Now in SQL with grouping

E(S,T)



- “Find nodes that have at least two distinct neighbors” (query)

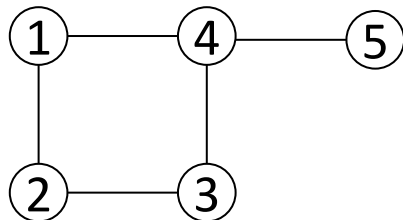


```
SELECT DISTINCT S
FROM E
GROUP BY S
HAVING COUNT(T) >= 2
```

```
SELECT DISTINCT E1.S
FROM E E1, E E2
WHERE E1.S = E2.S
AND E1.T <> E2.T
```

{1, 2, 3, 4}

- “each node has at least two distinct neighbors”



```
SELECT not exists
(SELECT *
FROM E E1
WHERE not exists
(SELECT *
FROM E E2
WHERE E1.S = E2.S
AND E1.T <> E2.T))
```

FALSE

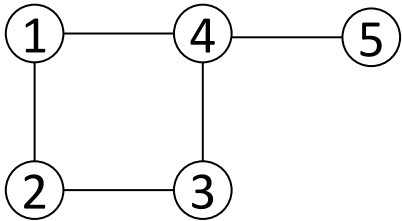
Now in SQL with grouping

E(S,T)



501

- “Find nodes that have at least two distinct neighbors” (query)

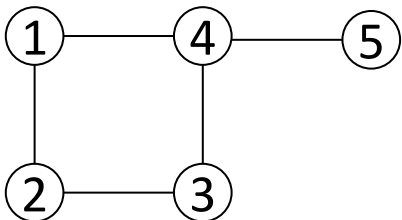


```
SELECT DISTINCT S
FROM E
GROUP BY S
HAVING COUNT(T) >= 2
```

```
SELECT DISTINCT E1.S
FROM E E1, E E2
WHERE E1.S = E2.S
AND E1.T <> E2.T
```

{1, 2, 3, 4}

- “each node has at least two distinct neighbors”



```
SELECT not exists
(SELECT S
FROM E
GROUP BY S
HAVING COUNT(T)=1)
```

```
SELECT not exists
(SELECT *
FROM E E1
WHERE not exists
(SELECT *
FROM E E2
WHERE E1.S = E2.S
AND E1.T <> E2.T))
```

FALSE

More practice



- "A small, happy dog is at home"

?

- "Every small dog that is at home is happy."

?

- "Jiahui owns a small, happy dog"

?

- "Jiahui owns every small, happy dog."

?

One more example



- "There are infinitely many prime numbers"



Components of FOL: (1) Syntax = First-order language

- **Alphabet**: symbols in use

- Variables, constants, **function** symbols, **predicate** symbols, connectives, quantifiers, punctuation symbols

vocabulary

x Alice MotherOf(x) Fibrous(x,y) \wedge \exists ()

terms *relation b/w objects*

- **Term**: expression that stands for an element or object

- Variable, constant
- Inductively $f(t_1, \dots, t_n)$ where t_i are terms, f a function symbol

MotherOf(MotherOf(x))

- **(Well-formed) formula**: parameterized statement

- **Atom** $p(t_1, \dots, t_n)$ where p is a predicate symbol, t_i terms (**atomic formula**, together with predicates $t_1=t_2$)
- Inductively, for formulas F, G , variable x :

$x = \text{'Alice'}$

$F \wedge G$ $F \vee G$ $\neg F$ $F \rightarrow G$ $F \leftrightarrow G$ $\forall x F$ $\exists x F$

- A first-order **language** refers to the set of all formulas over an alphabet