# CSE 544: Principles of Database Systems

**Query Complexity** 

#### Announcements

• Project milestone due tomorrow night

• HW3 due on Monday

# Query Complexity

Fix some query language L

 What is the complexity of evaluating queries in L on input databases?

- The answer gives us several insights:
  - What physical operators we need for L
  - What queries can and cannot we write in  $\ensuremath{\mathsf{L}}$
  - Whether we need to extend L

#### We are interested in these classes



# **Measuring Size**

 The size of a database D is the number of tuples



<u> </u>	
С	b
а	С
а	а
b	а
b	b
g	b
f	а

S =

 Alternatively, the size of D is the number of constants in its active domain

Number of tuples: 5+7 = 12

ADom(D)={a,b,c,d,e,f,g} Size of the active domain = 7

#### These two are related (how?)



Give an algorithm for computing Q on any input **D**. Express its complexity as a function of n = |ADom(D)|

\* Active Domain = all constants in D. ADom(D) = {a,b,c,...}



endfor



endfor

![](_page_8_Figure_1.jpeg)

```
for x in ADom do
  good_v = false
  for y in ADom do
    if (x,y) \in \mathbb{R} then
        good_7 = true
        for z in ADom do
           if (y,z) \in S then
               good_u = false
               for u in Adom do
                  if (z,u) \in \mathbb{R} then good, = true
               endfor
               if not good, then good_7 = false
         endfor
         if good_z then good_v = true
   endfor
   if good<sub>v</sub> then output x
endfor
```

What is the running time f(n)?

4=number of variables

In what complexity class is this query?

![](_page_8_Picture_7.jpeg)

![](_page_9_Figure_1.jpeg)

```
for x in ADom do
  good_v = false
  for y in ADom do
    if (x,y) \in \mathbb{R} then
        good_7 = true
        for z in ADom do
           if (y,z) \in S then
               good_u = false
               for u in Adom do
                 if (z,u) \in \mathbb{R} then good, = true
               endfor
               if not good, then good_7 = false
        endfor
        if good_z then good_v = true
   endfor
  if good_v then output x
endfor
```

What is the running time f(n)?

4=number of variables

In what complexity class is this query?

![](_page_9_Picture_7.jpeg)

![](_page_10_Figure_1.jpeg)

```
for x in ADom do
  good_v = false
  for y in ADom do
    if (x,y) \in \mathbb{R} then
        good_7 = true
       for z in ADom do
           if (y,z) \in S then
               good_u = false
               for u in Adom do
                 if (z,u) \in \mathbb{R} then good, = true
               endfor
               if not good, then good_7 = false
        endfor
        if good_z then good_v = true
   endfor
  if good_v then output x
endfor
```

What is the running time f(n)?

4=number of variables

In what complexity class is this query?

![](_page_10_Picture_7.jpeg)

### Discussion

- The query complexity of O(n<sup>4</sup>) if we assumed that the query is fixed and only the database is variable
- If we assume that both query and database are variable, then the complexity is different (and much higher)

• Data complexity v.s. Query complexity

Vardi, *The Complexity of Relational Query Languages*, STOC 1982

Query Q, database D

- <u>Data complexity</u>: fix Q, complexity = f(D)
- <u>Query complexity</u>: fix D, complexity = f(Q)
- <u>Combined complexity</u>: complexity = f(D,Q)

![](_page_12_Picture_5.jpeg)

Moshe Vardi 2008 ACM SIGMOD Codd Innovation Award

# Data Complexity

Given a query Q in a query language L what is the complexity of the following problem? "Given D, compute Q(D)"

- Language design tradeoff
  - High complexity  $\rightarrow$  L can express rich queries
  - Low complexity  $\rightarrow$  L can be implemented efficiently

# Conventions

- The complexity is usually defined for a decision problem
  - Hence, we will study only the complexity of Boolean queries
- The complexity usually assumes some encoding of the input
  - Hence we will encode the database instances using a binary representation

**Definition** A Boolean Query is a query that returns either true or false

**<u>Factoid</u>**: the complexity of <u>L</u> is fully determined by its Boolean queries

**Definition** A Boolean Query is a query that returns either true or false

**Factoid**: the complexity of L is fully determined by its Boolean queries

#### Non-boolean queries

Boolean queries:

$$Q = \exists z.R(a',z) \land S(z,b') \quad a,b \in ADom(D)$$

 $Q(x,y) = \exists z.R(x,z) \land S(z,y)$ 

**Definition** A Boolean Query is a query that returns either true or false

**Factoid**: the complexity of L is fully determined by its Boolean queries

#### Non-boolean queries

 $Q(x,y) = \exists z.R(x,z) \land S(z,y)$ 

SELECT DISTINCT R.x, S.y FROM R, S WHERE R.z = S.z Boolean queries:

$$Q = \exists z.R(a',z) \land S(z,b')$$

a,b ∈ ADom(D)

SELECT DISTINCT 'yes' FROM R, S WHERE R.x='a' and R.y = S.y and S.y='b'

**Definition** A Boolean Query is a query that returns either true or false

**Factoid**: the complexity of L is fully determined by its Boolean queries

![](_page_18_Figure_3.jpeg)

 $Q(x,y) = \exists z.R(x,z) \land S(z,y)$ 

SELECT DISTINCT R.x, S.y FROM R, S WHERE R.z = S.z

 $Q(x)=\exists y.R(x,y) \land (\forall z.S(y,z) \rightarrow \exists u.R(z,u))$ 

Boolean queries:

$$Q = \exists z.R(a',z) \land S(z,b')$$

 $a,b \in ADom(D)$ 

```
SELECT DISTINCT 'yes'
FROM R, S
WHERE R.x='a' and R.y = S.y and S.y='b'
```

 $Q=\exists y.R(a',y) \land (\forall z.S(y,z) \rightarrow \exists u.R(z,u))$ 

**Definition** A Boolean Query is a query that returns either true or false

**<u>Factoid</u>**: the complexity of <u>L</u> is fully determined by its Boolean queries

Non-boolean queries	Boolean queries:
$Q(x,y) = \exists z.R(x,z) \land S(z,y)$	$Q = \exists z.R(a',z) \land S(z,b') a,b \in ADom(D)$
SELECT DISTINCT R.x, S.y FROM R, S WHERE R.z = S.z	SELECT DISTINCT 'yes' FROM R, S WHERE R.x='a' and R.y = S.y and S.y='b'
$Q(x)=\exists y.R(x,y) \land (\forall z.S(y,z) \rightarrow \exists u.R(z,u))$	$Q=\exists y.R(a',y) \land (\forall z.S(y,z) \rightarrow \exists u.R(z,u))$
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$

# **Database Encoding**

 We could encode a database D as a list of encodings of tuples, similar to how database systems store data on disk: – "[R(a,b),R(cac,bbfg),R(cac,b)],[S(b),S(cac)]"

• Instead, we use a simpler representation, generalizing adjacency matrix for graphs

### **Database Encoding**

Encode **D** = (D,  $R_1^D$ , ...,  $R_k^D$ ) as follows:

- Let n = |ADom(D)|
- If R<sub>i</sub> has arity k, then encode it as a string of n<sup>k</sup> bits:
  - -0 means element  $(a_1, \ldots, a_k) \notin R_i^D$
  - 1 means element  $(a_1, ..., a_k) \in R_i^D$

![](_page_21_Figure_6.jpeg)

![](_page_21_Figure_7.jpeg)

![](_page_21_Figure_8.jpeg)

# The Data Complexity

Fix any Boolean query Q in the query language. Determine the complexity of the following problem:

- Given an input database instance D = (D, R<sub>1</sub><sup>D</sup>, ..., R<sub>k</sub><sup>D</sup>), check if Q(D) = true.
- This is also known as the <u>Model Checking</u> <u>Problem</u>: check if **D** is a model for Q.

# What we will discuss next time

- Relational queries
- Datalog and stratified datalog<sup>¬</sup>
- Datalog<sup>¬</sup> with inflationary fixpoint
- Datalog<sup>¬</sup> with partial fixpoint

- AC<sup>0</sup>
- L = a.k.a. LOGSPACE
- NL= a.k.a. NLOGSPACE
- NC
- P = a.k.a. PTIME
- NP
- PSPACE