

# Finite Model Theory

## Lecture 4: Conjunctive Queries

Spring 2025

# Where We Are

- Thanks for submitting HW1.
- HW2 still to be released: it covers this lecture and the next one.
- Demystifying **decidability** of  $\Sigma$ . Two problems:
  - Membership: given  $\varphi$ , is  $\varphi \in \Sigma$ ?
  - Implication: given  $\varphi$ , does  $\Sigma \models \varphi$  hold?
  - If membership is r.e. then implication is r.e. (Gödel's **completeness**).
  - Example of  $\Sigma$  that is not r.e.:  $\text{Th}(\mathbb{N}, +, *)$  (Gödel's **incompleteness**).

**Today:** Model Checking, and Conjunctive Queries.

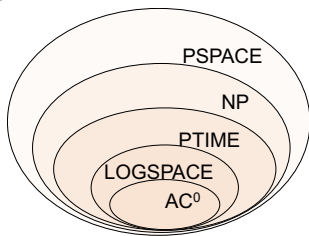
# Model Checking

# Model Checking

Model checking: given  $\mathbf{A}$ ,  $\varphi$ , check whether  $\mathbf{A} \models \varphi$ .

When  $\mathbf{A}$  is a database and  $\varphi$  a query, then this is called *query evaluation*. This is the bread-and-butter of database engines.

What is the complexity of checking  $\mathbf{A} \models \varphi$ ?



# Three Ways to Define Complexity

$$\mathbf{A} \models \varphi$$

- **Data Complexity**. Fix the sentence  $\varphi$ , complexity is  $f(|\mathbf{A}|)$ .
- **Expression Complexity**. Fix the structure  $\mathbf{A}$ , complexity is  $f(|\varphi|)$ .
- **Combined Complexity**,  $f(|\varphi|, |\mathbf{A}|)$ .

# Three Ways to Define Complexity

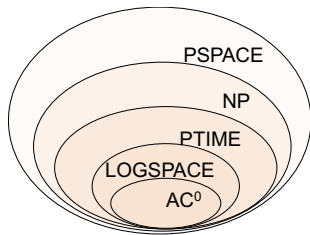
$$\mathbf{A} \models \varphi$$

## Theorem

The *Data* Complexity of FO is in  $AC^0$  (actually in uniform  $AC^0$ ).

The *Expression- and Combined* Complexity of FO is PSPACE complete.

$AC^0$  is at the bottom of the hierarchy  
We start by proving something simpler:  
The Data Complexity of FO is in PTIME.



## Data Complexity of FO is in PTIME: Proof

How do we evaluate this?  $\varphi = \exists x (A(x) \wedge \forall y (B(y) \Rightarrow C(x, y)))$

Assume the domain is  $[N]$ .

# Data Complexity of FO is in PTIME: Proof

How do we evaluate this?  $\varphi = \exists x (A(x) \wedge \forall y (B(y) \Rightarrow C(x, y)))$

Assume the domain is  $[N]$ .

```
fx = false;
for x = 1, N do:
  if A(x) then:
    fy = true
    for y = 1, N do:
      if not (B(y) => C(x,y))
        then: fy = false;
    if fy then: fx = true;
return fx
```



# Data Complexity of FO is in PTIME: Proof

How do we evaluate this?  $\varphi = \exists x (A(x) \wedge \forall y (B(y) \Rightarrow C(x, y)))$

Assume the domain is  $[N]$ .

```
fx = false;
for x = 1, N do:
  if A(x) then:
    fy = true
    for y = 1, N do:
      if not (B(y) => C(x,y))
        then: fy = false;
    if fy then: fx = true;
return fx
```

- Generalizes to any sentence in prenex normal form  $\varphi$
- Runtime  $O(N^k)$ , where:  
 $N = |\text{Dom}|$ ,  $k = |\text{Vars}(\varphi)|$
- In PTIME (and in LOGSPACE).

# Data Complexity of FO is in PTIME: Proof

How do we evaluate this?  $\varphi = \exists x (A(x) \wedge \forall y (B(y) \Rightarrow C(x, y)))$

Assume the domain is  $[N]$ .

```
fx = false;
for x = 1, N do:
  if A(x) then:
    fy = true
    for y = 1, N do:
      if not (B(y) => C(x,y))
        then: fy = false;
    if fy then: fx = true;
return fx
```

- Generalizes to any sentence in prenex normal form  $\varphi$
- Runtime  $O(N^k)$ , where:  
 $N = |\text{Dom}|$ ,  $k = |\text{Vars}(\varphi)|$
- In PTIME (and in LOGSPACE).

Many texts state that the data complexity is in LOGSPACE, or in PTIME. The correct complexity is  $AC^0$ . **Let's prove it**

# Definition of $AC^0$

## Definition

A problem is in  $AC^0$  if  $\forall N$ , there exists a circuit of polynomial size and constant depth, consisting NOT gates and unbounded fan-in AND and OR gates, that computes the problem when the input is encoded using  $N$  bits.

# Definition of $AC^0$

## Definition

A problem is in  $AC^0$  if  $\forall N$ , there exists a circuit of polynomial size and constant depth, consisting NOT gates and unbounded fan-in AND and OR gates, that computes the problem when the input is encoded using  $N$  bits.

E.g. Given an undirected graph with  $n$  nodes, check if it has a triangle.

# Definition of $AC^0$

## Definition

A problem is in  $AC^0$  if  $\forall N$ , there exists a circuit of polynomial size and constant depth, consisting NOT gates and unbounded fan-in AND and OR gates, that computes the problem when the input is encoded using  $N$  bits.

E.g. Given an undirected graph with  $n$  nodes, check if it has a triangle.  
When  $n = 4$  there are  $N = 6$  possible edges,  $E_{ij}$  for  $1 \leq i < j \leq 4$

# Definition of $AC^0$

## Definition

A problem is in  $AC^0$  if  $\forall N$ , there exists a circuit of polynomial size and constant depth, consisting NOT gates and unbounded fan-in AND and OR gates, that computes the problem when the input is encoded using  $N$  bits.

E.g. Given an undirected graph with  $n$  nodes, check if it has a triangle.  
When  $n = 4$  there are  $N = 6$  possible edges,  $E_{ij}$  for  $1 \leq i < j \leq 4$

$$(E_{12} \wedge E_{23} \wedge E_{13}) \vee (E_{12} \wedge E_{24} \wedge E_{14}) \vee (E_{23} \wedge E_{34} \wedge E_{24}) \vee (E_{34} \wedge E_{14} \wedge E_{13})$$

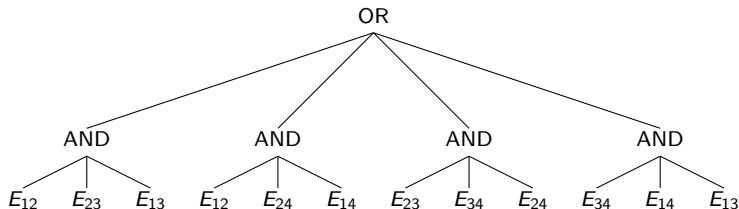
# Definition of $AC^0$

## Definition

A problem is in  $AC^0$  if  $\forall N$ , there exists a circuit of polynomial size and constant depth, consisting NOT gates and unbounded fan-in AND and OR gates, that computes the problem when the input is encoded using  $N$  bits.

E.g. Given an undirected graph with  $n$  nodes, check if it has a triangle.  
When  $n = 4$  there are  $N = 6$  possible edges,  $E_{ij}$  for  $1 \leq i < j \leq 4$

$$(E_{12} \wedge E_{23} \wedge E_{13}) \vee (E_{12} \wedge E_{24} \wedge E_{14}) \vee (E_{23} \wedge E_{34} \wedge E_{24}) \vee (E_{34} \wedge E_{14} \wedge E_{13})$$



# Data Complexity FO is in $AC^0$ : Proof

Fix  $\varphi$  in FO. Encode the input  $\mathbf{A}$  using bits:

- Let  $N = |\text{Dom}(\mathbf{A})|$ .



## Data Complexity FO is in $AC^0$ : Proof

Fix  $\varphi$  in FO. Encode the input  $\mathbf{A}$  using bits:

- Let  $N = |\text{Dom}(\mathbf{A})|$ .
- Encode a relation of arity  $k$  using  $N^k$  bits.  
E.g.  $R(X, Y)$  encoded using Boolean matrix  $R_{ij}$ , with  $N^2$  bits.

# Data Complexity FO is in $AC^0$ : Proof

Fix  $\varphi$  in FO. Encode the input  $\mathbf{A}$  using bits:

- Let  $N = |\text{Dom}(\mathbf{A})|$ .
- Encode a relation of arity  $k$  using  $N^k$  bits.  
E.g.  $R(X, Y)$  encoded using Boolean matrix  $R_{ij}$ , with  $N^2$  bits.
- Expand  $\varphi$  into a Boolean formula  $F_{\varphi, \mathbf{A}}$ , called the **lineage** of  $\varphi$  on  $\mathbf{A}$ .

# Data Complexity FO is in $AC^0$ : Proof

Fix  $\varphi$  in FO. Encode the input  $\mathbf{A}$  using bits:

- Let  $N = |\text{Dom}(\mathbf{A})|$ .
- Encode a relation of arity  $k$  using  $N^k$  bits.  
E.g.  $R(X, Y)$  encoded using Boolean matrix  $R_{ij}$ , with  $N^2$  bits.
- Expand  $\varphi$  into a Boolean formula  $F_{\varphi, \mathbf{A}}$ , called the **lineage** of  $\varphi$  on  $\mathbf{A}$ .
- Represent  $F_{\varphi, \mathbf{A}}$  using an  $AC^0$  circuit.

## Data Complexity FO is in $AC^0$ : Proof

Fix  $\varphi$  in FO. Encode the input  $\mathbf{A}$  using bits:

- Let  $N = |\text{Dom}(\mathbf{A})|$ .
- Encode a relation of arity  $k$  using  $N^k$  bits.  
E.g.  $R(X, Y)$  encoded using Boolean matrix  $R_{ij}$ , with  $N^2$  bits.
- Expand  $\varphi$  into a Boolean formula  $F_{\varphi, \mathbf{A}}$ , called the **lineage** of  $\varphi$  on  $\mathbf{A}$ .
- Represent  $F_{\varphi, \mathbf{A}}$  using an  $AC^0$  circuit.

Example:  $\varphi = \exists x (A(x) \wedge \forall y (B(y) \Rightarrow C(x, y)))$

## Data Complexity FO is in $AC^0$ : Proof

Fix  $\varphi$  in FO. Encode the input  $\mathbf{A}$  using bits:

- Let  $N = |\text{Dom}(\mathbf{A})|$ .
- Encode a relation of arity  $k$  using  $N^k$  bits.  
E.g.  $R(X, Y)$  encoded using Boolean matrix  $R_{ij}$ , with  $N^2$  bits.
- Expand  $\varphi$  into a Boolean formula  $F_{\varphi, \mathbf{A}}$ , called the **lineage** of  $\varphi$  on  $\mathbf{A}$ .
- Represent  $F_{\varphi, \mathbf{A}}$  using an  $AC^0$  circuit.

Example:  $\varphi = \exists x (A(x) \wedge \forall y (B(y) \Rightarrow C(x, y)))$

Its **lineage** is:  $F_{\varphi, \mathbf{A}} = \bigvee_{i=1, N} (A_i \wedge \bigwedge_{j=1, N} (\neg B_j \vee C_{ij}))$

## Data Complexity FO is in $AC^0$ : Proof

Fix  $\varphi$  in FO. Encode the input  $\mathbf{A}$  using bits:

- Let  $N = |\text{Dom}(\mathbf{A})|$ .
- Encode a relation of arity  $k$  using  $N^k$  bits.  
E.g.  $R(X, Y)$  encoded using Boolean matrix  $R_{ij}$ , with  $N^2$  bits.
- Expand  $\varphi$  into a Boolean formula  $F_{\varphi, \mathbf{A}}$ , called the **lineage** of  $\varphi$  on  $\mathbf{A}$ .
- Represent  $F_{\varphi, \mathbf{A}}$  using an  $AC^0$  circuit.

Example:  $\varphi = \exists x (A(x) \wedge \forall y (B(y) \Rightarrow C(x, y)))$

Its **lineage** is:  $F_{\varphi, \mathbf{A}} = \bigvee_{i=1, N} (A_i \wedge \bigwedge_{j=1, N} (\neg B_j \vee C_{ij}))$

Represented by a circuit of depth 5 and size  $O(N^2)$  **in class**

# Summary

- **Data complexity:**  $AC^0$ . This implies LOGSPACE, PTIME.
  - $AC^0$  is the class of highly parallelizable problems.
  - “SQL is embarrassingly parallel”
- **Expression complexity, combined complexity:** PSPACE complete
  - Membership is straightforward **WHY????**
  - We will prove hardness later.

# Conjunctive Queries



# Motivation

- FO is too rich for static analysis: Trakhtenbrot's theorem is a fundamental limit.
- For some FO fragments static analysis is possible, and they still capture the most important queries in practice.
- Conjunctive Queries (CQ) and Unions of Conjunctive Queries (UCQ).

# Motivation

- FO is too rich for static analysis: Trakhtenbrot's theorem is a fundamental limit.
- For some FO fragments static analysis is possible, and they still capture the most important queries in practice.
- Conjunctive Queries (CQ) and Unions of Conjunctive Queries (UCQ).

Terminology  
alert

query=formula;  
Boolean query=sentence;

database instance=structure;  
query evaluation=model checking

# Conjunctive Queries

## Definition

A **Conjunctive Query (CQ)** is a formula of the form:

$$Q(\mathbf{x}_0) = \exists \mathbf{y} (R_1(\mathbf{x}_1) \wedge R_2(\mathbf{x}_2) \wedge \cdots \wedge R_m(\mathbf{x}_m))$$

---

<sup>1</sup>Relational Algebra:  $\sigma, \Pi, \bowtie, \cup, -$ .

# Conjunctive Queries

## Definition

A **Conjunctive Query (CQ)** is a formula of the form:

$$Q(\mathbf{x}_0) = \exists \mathbf{y} (R_1(\mathbf{x}_1) \wedge R_2(\mathbf{x}_2) \wedge \cdots \wedge R_m(\mathbf{x}_m))$$

- The free variables,  $\mathbf{x}_0$ , are called **head variables**.

---

<sup>1</sup>Relational Algebra:  $\sigma, \Pi, \bowtie, \cup, -$ .

# Conjunctive Queries

## Definition

A **Conjunctive Query (CQ)** is a formula of the form:

$$Q(\mathbf{x}_0) = \exists \mathbf{y} (R_1(\mathbf{x}_1) \wedge R_2(\mathbf{x}_2) \wedge \cdots \wedge R_m(\mathbf{x}_m))$$

- The free variables,  $\mathbf{x}_0$ , are called **head variables**.
- **Safety Condition**: every head variable must occur in some atom  $R_j$ .

---

<sup>1</sup>Relational Algebra:  $\sigma, \Pi, \bowtie, \cup, -$ .

# Conjunctive Queries

## Definition

A **Conjunctive Query (CQ)** is a formula of the form:

$$Q(\mathbf{x}_0) = \exists \mathbf{y} (R_1(\mathbf{x}_1) \wedge R_2(\mathbf{x}_2) \wedge \cdots \wedge R_m(\mathbf{x}_m))$$

- The free variables,  $\mathbf{x}_0$ , are called **head variables**.
- **Safety Condition**: every head variable must occur in some atom  $R_j$ .
- E.g.  $Q(x, y) = \exists z (E(x, z) \wedge E(z, y))$ .

---

<sup>1</sup>Relational Algebra:  $\sigma, \Pi, \bowtie, \cup, -$ .

# Conjunctive Queries

## Definition

A **Conjunctive Query (CQ)** is a formula of the form:

$$Q(\mathbf{x}_0) = \exists \mathbf{y} (R_1(\mathbf{x}_1) \wedge R_2(\mathbf{x}_2) \wedge \cdots \wedge R_m(\mathbf{x}_m))$$

- The free variables,  $\mathbf{x}_0$ , are called **head variables**.
- **Safety Condition**: every head variable must occur in some atom  $R_j$ .
- E.g.  $Q(x, y) = \exists z (E(x, z) \wedge E(z, y))$ .
- Same as FO formulas restricted to  $=, \wedge, \exists$   
Same as RA<sup>1</sup> restricted to  $\sigma, \Pi, \bowtie$ .  
Same as SELECT-DISTINCT-FROM-WHERE

---

<sup>1</sup>Relational Algebra:  $\sigma, \Pi, \bowtie, \cup, -$ .

# Unions of Conjunctive Queries

## Definition

A **Union of Conjunctive Queries (UCQ)** is a formula of the form:

$$Q(\mathbf{x}) = Q_1(\mathbf{x}) \vee Q_2(\mathbf{x}) \vee \cdots \vee Q_k(\mathbf{x})$$

where all  $Q_i$ 's are CQs, and have the same sets of free variables.



# Unions of Conjunctive Queries

## Definition

A **Union of Conjunctive Queries (UCQ)** is a formula of the form:

$$Q(\mathbf{x}) = Q_1(\mathbf{x}) \vee Q_2(\mathbf{x}) \vee \dots \vee Q_k(\mathbf{x})$$

where all  $Q_i$ 's are CQs, and have the same sets of free variables.

- E.g.  $Q(x, y) = E(x, y) \vee \exists z(E(x, z) \wedge E(z, y))$ .

# Unions of Conjunctive Queries

## Definition

A **Union of Conjunctive Queries (UCQ)** is a formula of the form:

$$Q(\mathbf{x}) = Q_1(\mathbf{x}) \vee Q_2(\mathbf{x}) \vee \dots \vee Q_k(\mathbf{x})$$

where all  $Q_i$ 's are CQs, and have the same sets of free variables.

- E.g.  $Q(x, y) = E(x, y) \vee \exists z(E(x, z) \wedge E(z, y))$ .
- Same as FO formulas restricted to  $=, \wedge, \exists, \vee$ .  
Same as RA restricted to  $\sigma, \Pi, \bowtie, \cup$ .

# Unions of Conjunctive Queries

## Definition

A **Union of Conjunctive Queries (UCQ)** is a formula of the form:

$$Q(\mathbf{x}) = Q_1(\mathbf{x}) \vee Q_2(\mathbf{x}) \vee \cdots \vee Q_k(\mathbf{x})$$

where all  $Q_i$ 's are CQs, and have the same sets of free variables.

- E.g.  $Q(x, y) = E(x, y) \vee \exists z(E(x, z) \wedge E(z, y))$ .
- Same as FO formulas restricted to  $=, \wedge, \exists, \vee$ .  
Same as RA restricted to  $\sigma, \Pi, \bowtie, \cup$ .
- Notice: all queries  $Q_i(\mathbf{x})$  must have the same head variables  $\mathbf{x}$ .

## Monotone Queries

Given two databases  $D, D'$  over the same schema, we write  $D \subseteq D'$  if  $R_i^D \subseteq R_i^{D'}$  for every relation  $R_i$  in the schema.

### Definition

A query  $Q$  is **monotone** if  $D \subseteq D'$  implies  $Q(D) \subseteq Q(D')$ .

**Example:**  $\exists x, y, z (E(x, y) \wedge E(y, z))$     **Non-example:**  $\exists x \forall y (V(x) \wedge \forall y (V(y) \Rightarrow E(x, y)))$

## Monotone Queries

Given two databases  $D, D'$  over the same schema, we write  $D \subseteq D'$  if  $R_i^D \subseteq R_i^{D'}$  for every relation  $R_i$  in the schema.

### Definition

A query  $Q$  is **monotone** if  $D \subseteq D'$  implies  $Q(D) \subseteq Q(D')$ .

**Example:**  $\exists x, y, z (E(x, y) \wedge E(y, z))$     **Non-example:**  $\exists x \forall y (V(y) \Rightarrow E(x, y))$

All UCQ queries are monotone. **Exercise**

The only non-monotone operators are:

- negation  $\neg$  in FO.
- difference  $-$  in RA.

# Terminology

- **Boolean query**: no head vars:  $Q() = \exists x \exists y \exists z (E(x, y) \wedge E(y, z)).$
- **Full query**: no existential vars:  $Q(x, y, z) = E(x, y) \wedge E(y, z).$
- **Without selfjoins**: every relation name occurs at most once.

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

- We often omit the existential quantifiers, and write for example:

$$Q(x) = R(x, y) \wedge S(y, z) \wedge T(z, x).$$

## Discussion

- A Boolean query  $Q$  is a sentence.

“Evaluating  $Q$  on a database  $\mathbf{D}$ ” means checking  $\mathbf{D} \models Q$ .

We often write  $Q(\mathbf{D}) = 1$  or  $Q(\mathbf{D}) = 0$  respectively.

## Discussion

- A Boolean query  $Q$  is a sentence.

“Evaluating  $Q$  on a database  $\mathbf{D}$ ” means checking  $\mathbf{D} \models Q$ .

We often write  $Q(\mathbf{D}) = 1$  or  $Q(\mathbf{D}) = 0$  respectively.

- A query  $Q(\mathbf{x})$  with head variables is a formula.

The “answer of  $Q$  on  $\mathbf{D}$ ” is  $\{\mathbf{a} \mid \mathbf{a} \in (\text{Dom}(\mathbf{D}))^{|\mathbf{x}|}, \mathbf{D} \models Q[\mathbf{a}/\mathbf{x}]\}$ .

We write  $Q(\mathbf{D})$  for the set of answers,  $Q(\mathbf{D}) \subseteq |\text{Dom}(\mathbf{D})|^{|\mathbf{x}|}$ .



# Model Checking for CQ

# Motivation

We already know that the data complexity is in  $AC^0$ .

What is the expression complexity? The combined complexity?

Will answer both, and also discuss the expression/combined complexity for FO (which we left out).

Importantly: we will define query evaluation for CQ in terms of  
**Homomorphisms**

## Equivalent Concepts

- A Conjunctive Query:

$$R(x, y, z) \wedge S(x, u) \wedge S(y, v) \wedge S(z, w) \wedge R(u, v, w)$$

## Equivalent Concepts

- A Conjunctive Query:

$$R(x, y, z) \wedge S(x, u) \wedge S(y, v) \wedge S(z, w) \wedge R(u, v, w)$$

- A structure  $\mathbf{A}$ , a.k.a. a database instance:

$$\text{Dom}(\mathbf{A}) = \{x, y, z, u, v, w\} \quad R^{\mathbf{A}} = \begin{array}{|c|c|c|} \hline x & y & z \\ \hline u & v & w \\ \hline \end{array} \quad S^{\mathbf{A}} = \begin{array}{|c|c|} \hline x & u \\ \hline y & v \\ \hline z & w \\ \hline \end{array}$$

## Equivalent Concepts

- A Conjunctive Query:

$$R(x, y, z) \wedge S(x, u) \wedge S(y, v) \wedge S(z, w) \wedge R(u, v, w)$$

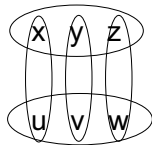
- A structure  $\mathbf{A}$ , a.k.a. a database instance:

$$\text{Dom}(\mathbf{A}) = \{x, y, z, u, v, w\} \quad R^{\mathbf{A}} = \begin{array}{|c|c|c|} \hline x & y & z \\ \hline u & v & w \\ \hline \end{array} \quad S^{\mathbf{A}} = \begin{array}{|c|c|} \hline x & u \\ \hline y & v \\ \hline z & w \\ \hline \end{array}$$

- A labeled hypergraph,  $G = (V, E)$ , where

$$V = \{x, y, z, u, v, w\}, \quad E = \{\{x, y, z\}, \{u, v, w\}, \{x, u\}, \{y, v\}, \{z, w\}\}$$

(hyperedges are labeled with  $R, S$  respectively).



## Equivalent Concepts

- A Conjunctive Query:

$$R(x, y, z) \wedge S(x, u) \wedge S(y, v) \wedge S(z, w) \wedge R(u, v, w)$$

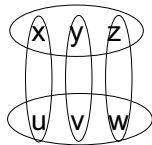
- A structure  $\mathbf{A}$ , a.k.a. a database instance:

$$\text{Dom}(\mathbf{A}) = \{x, y, z, u, v, w\} \quad R^{\mathbf{A}} = \begin{array}{|c|c|c|} \hline x & y & z \\ \hline u & v & w \\ \hline \end{array} \quad S^{\mathbf{A}} = \begin{array}{|c|c|} \hline x & u \\ \hline y & v \\ \hline z & w \\ \hline \end{array}$$

- A labeled hypergraph,  $G = (V, E)$ , where

$$V = \{x, y, z, u, v, w\}, \quad E = \{\{x, y, z\}, \{u, v, w\}, \{x, u\}, \{y, v\}, \{z, w\}\}$$

(hyperedges are labeled with  $R, S$  respectively).



We will often switch back-and-forth between these equivalent notions

# Homomorphisms

$$\mathbf{A} = (\text{Dom}(A), R_1^A, \dots, R_m^A); \mathbf{B} = (\text{Dom}(B), R_1^B, \dots, R_m^B).$$

## Definition

A **homomorphism**  $h : \mathbf{A} \rightarrow \mathbf{B}$  is a function  $h : \text{Dom}(A) \rightarrow \text{Dom}(B)$  such that, for all  $j = 1, m$ ,  $h(R_j^A) \subseteq R_j^B$ .

# Homomorphisms

$\mathbf{A} = (\text{Dom}(A), R_1^A, \dots, R_m^A); \mathbf{B} = (\text{Dom}(B), R_1^B, \dots, R_m^B).$

## Definition

A **homomorphism**  $h : \mathbf{A} \rightarrow \mathbf{B}$  is a function  $h : \text{Dom}(A) \rightarrow \text{Dom}(B)$  such that, for all  $j = 1, m$ ,  $h(R_j^A) \subseteq R_j^B$ .

Graph homomorphism  $h : G \rightarrow G'$  is  $h : V \rightarrow V'$  s.t.  $\forall e \in E, h(e) \in E'$ .



# Homomorphisms

$\mathbf{A} = (\text{Dom}(A), R_1^A, \dots, R_m^A)$ ;  $\mathbf{B} = (\text{Dom}(B), R_1^B, \dots, R_m^B)$ .

## Definition

A **homomorphism**  $h : \mathbf{A} \rightarrow \mathbf{B}$  is a function  $h : \text{Dom}(A) \rightarrow \text{Dom}(B)$  such that, for all  $j = 1, m$ ,  $h(R_j^A) \subseteq R_j^B$ .

Graph homomorphism  $h : G \rightarrow G'$  is  $h : V \rightarrow V'$  s.t.  $\forall e \in E, h(e) \in E'$ .

For query homomorphisms  $h : Q(\mathbf{x}) \rightarrow Q'(\mathbf{x}')$  we also require  $h$  to map head variables to head variables and constants to constants.

## Query Evaluation for CQ and Homomorphisms

If  $Q$  is Boolean, then  $\mathbf{D} \models Q$  iff  $\exists h : Q \rightarrow \mathbf{D}$ .

If  $Q(\mathbf{x})$  has head variables, then  $Q(\mathbf{D}) = \{h(\mathbf{x}) \mid h : Q \rightarrow \mathbf{D}\}$ .

## Query Evaluation for CQ and Homomorphisms

If  $Q$  is Boolean, then  $\mathbf{D} \models Q$  iff  $\exists h : Q \rightarrow \mathbf{D}$ .

If  $Q(\mathbf{x})$  has head variables, then  $Q(\mathbf{D}) = \{h(\mathbf{x}) \mid h : Q \rightarrow \mathbf{D}\}$ .

$$Q(x) = R(x) \wedge S(x, y) \wedge T(y, 'a')$$

# Query Evaluation for CQ and Homomorphisms

If  $Q$  is Boolean, then  $\mathbf{D} \models Q$  iff  $\exists h : Q \rightarrow \mathbf{D}$ .

If  $Q(\mathbf{x})$  has head variables, then  $Q(\mathbf{D}) = \{h(\mathbf{x}) \mid h : Q \rightarrow \mathbf{D}\}$ .

$$Q(x) = R(x) \wedge S(x, y) \wedge T(y, 'a')$$

 $R =$ 

x
1
2
3

 $S =$ 

x	y
1	10
1	20
2	20

 $T =$ 

y	z
10	<i>a</i>
10	<i>b</i>
20	<i>a</i>

## Query Evaluation for CQ and Homomorphisms

If  $Q$  is Boolean, then  $\mathbf{D} \models Q$  iff  $\exists h : Q \rightarrow \mathbf{D}$ .

If  $Q(\mathbf{x})$  has head variables, then  $Q(\mathbf{D}) = \{h(\mathbf{x}) \mid h : Q \rightarrow \mathbf{D}\}$ .

$$Q(x) = R(x) \wedge S(x, y) \wedge T(y, 'a')$$

 $R =$ 

x
1
2
3

 $S =$ 

x	y
1	10
1	20
2	20

 $T =$ 

y	z
10	<i>a</i>
10	<i>b</i>
20	<i>a</i>

We list all homomorphisms  $h : \{x, y\} \rightarrow \{1, 2, 10, 20, 'a'\}$ :

# Query Evaluation for CQ and Homomorphisms

If  $Q$  is Boolean, then  $\mathbf{D} \models Q$  iff  $\exists h : Q \rightarrow \mathbf{D}$ .

If  $Q(\mathbf{x})$  has head variables, then  $Q(\mathbf{D}) = \{h(\mathbf{x}) \mid h : Q \rightarrow \mathbf{D}\}$ .

$$Q(x) = R(x) \wedge S(x, y) \wedge T(y, 'a')$$

$$R =$$

x
1
2
3

$$S =$$

x	y
1	10
1	20
2	20

$$T =$$

y	z
10	<i>a</i>
10	<i>b</i>
20	<i>a</i>

We list all homomorphisms  $h : \{x, y\} \rightarrow \{1, 2, 10, 20, 'a'\}$ :

$$h =$$

$x (= \text{Head}(Q))$	y	<i>a</i>
1	10	<i>a</i>

# Query Evaluation for CQ and Homomorphisms

If  $Q$  is Boolean, then  $\mathbf{D} \models Q$  iff  $\exists h : Q \rightarrow \mathbf{D}$ .

If  $Q(\mathbf{x})$  has head variables, then  $Q(\mathbf{D}) = \{h(\mathbf{x}) \mid h : Q \rightarrow \mathbf{D}\}$ .

$$Q(x) = R(x) \wedge S(x, y) \wedge T(y, 'a')$$

$$R =$$

x
1
2
3

$$S =$$

x	y
1	10
1	20
2	20

$$T =$$

y	z
10	<i>a</i>
10	<i>b</i>
20	<i>a</i>

We list all homomorphisms  $h : \{x, y\} \rightarrow \{1, 2, 10, 20, 'a'\}$ :

$$h =$$

$x(= \text{Head}(Q))$	y	<i>a</i>
1	10	<i>a</i>
1	20	<i>a</i>

# Query Evaluation for CQ and Homomorphisms

If  $Q$  is Boolean, then  $\mathbf{D} \models Q$  iff  $\exists h : Q \rightarrow \mathbf{D}$ .

If  $Q(\mathbf{x})$  has head variables, then  $Q(\mathbf{D}) = \{h(\mathbf{x}) \mid h : Q \rightarrow \mathbf{D}\}$ .

$$Q(x) = R(x) \wedge S(x, y) \wedge T(y, 'a')$$

$$R =$$

x
1
2
3

$$S =$$

x	y
1	10
1	20
2	20

$$T =$$

y	z
10	<i>a</i>
10	<i>b</i>
20	<i>a</i>

We list all homomorphisms  $h : \{x, y\} \rightarrow \{1, 2, 10, 20, 'a'\}$ :

$$h =$$

$x (= \text{Head}(Q))$	y	
1	10	<i>a</i>
1	20	<i>a</i>
2	20	<i>a</i>



# Query Evaluation for CQ and Homomorphisms

If  $Q$  is Boolean, then  $\mathbf{D} \models Q$  iff  $\exists h : Q \rightarrow \mathbf{D}$ .

If  $Q(\mathbf{x})$  has head variables, then  $Q(\mathbf{D}) = \{h(\mathbf{x}) \mid h : Q \rightarrow \mathbf{D}\}$ .

$$Q(x) = R(x) \wedge S(x, y) \wedge T(y, 'a')$$

$$R =$$

x
1
2
3

$$S =$$

x	y
1	10
1	20
2	20

$$T =$$

y	z
10	<i>a</i>
10	<i>b</i>
20	<i>a</i>

We list all homomorphisms  $h : \{x, y\} \rightarrow \{1, 2, 10, 20, 'a'\}$ :

$$h =$$

$x (= \text{Head}(Q))$	y	
1	10	<i>a</i>
1	20	<i>a</i>
2	20	<i>a</i>

Final answer after duplicate elimination:  $Q(\mathbf{D}) = \{1, 2\}$ .

# The Combined Complexity for UCQ is in NP

## Theorem

*The combined complexity of Boolean UCQ is in NP.*

**Proof:** Fix a UCQ  $Q = Q_1 \vee Q_2 \vee \dots$  and a database  $\mathbf{D}$ .

To check  $\mathbf{D} \models Q$ :

- “guess” a CQ  $Q_i$ , and
- “guess” a homomorphism  $h : Q_i \rightarrow \mathbf{D}$

# The Expression Complexity for CQ is NP-hard

## Theorem

*There exists a database  $\mathbf{D}$  for which the expression complexity of CQ queries is NP complete.*

Thus, both expression and combined complexities are NP-complete.

**Proof** Many proofs are possible (will explain shortly why). We will use reduction from 3SAT, because we will reuse it a few times.

# The Expression Complexity for CQ is NP-hard

## Theorem

*There exists a database  $\mathbf{D}$  for which the expression complexity of CQ queries is NP complete.*

Thus, both expression and combined complexities are NP-complete.

**Proof** Many proofs are possible (will explain shortly why). We will use reduction from 3SAT, because we will reuse it a few times.

Given a 3CNF formula  $\Phi$  we construct  $Q_\Phi, \mathbf{D}$  such that:

# The Expression Complexity for CQ is NP-hard

## Theorem

*There exists a database  $\mathbf{D}$  for which the expression complexity of CQ queries is NP complete.*

Thus, both expression and combined complexities are NP-complete.

**Proof** Many proofs are possible (will explain shortly why). We will use reduction from 3SAT, because we will reuse it a few times.

Given a 3CNF formula  $\Phi$  we construct  $Q_\Phi, \mathbf{D}$  such that:

$$\Phi \text{ is satisfiable iff } \exists h : Q_\Phi \rightarrow \mathbf{D}.$$

Notice that  $\mathbf{D}$  is independent of  $\Phi$ .

Details next.

## Reduction from 3SAT to CQ Evaluation

Given a 3CNF formula  $\Phi$  we construct  $Q_\Phi, \mathbf{D}$  such that:

$\Phi$  is satisfiable iff  $\exists h : Q_\Phi \rightarrow \mathbf{D}$ .

$Q_\Phi$  has one atom for each clause  $C$  in  $\Phi$ :

- If  $C = (X_i \vee X_j \vee X_k)$  then  $Q_\Phi$  contains  $A(x_i, x_j, x_k)$ .

## Reduction from 3SAT to CQ Evaluation

Given a 3CNF formula  $\Phi$  we construct  $Q_\Phi, \mathbf{D}$  such that:

$\Phi$  is satisfiable iff  $\exists h : Q_\Phi \rightarrow \mathbf{D}$ .

$Q_\Phi$  has one atom for each clause  $C$  in  $\Phi$ :

- If  $C = (X_i \vee X_j \vee X_k)$  then  $Q_\Phi$  contains  $A(x_i, x_j, x_k)$ .
- If  $C = (X_i \vee X_j \vee \neg X_k)$  then  $Q_\Phi$  contains  $B(x_i, x_j, x_k)$ .

## Reduction from 3SAT to CQ Evaluation

Given a 3CNF formula  $\Phi$  we construct  $Q_\Phi, \mathbf{D}$  such that:

$\Phi$  is satisfiable iff  $\exists h : Q_\Phi \rightarrow \mathbf{D}$ .

$Q_\Phi$  has one atom for each clause  $C$  in  $\Phi$ :

- If  $C = (X_i \vee X_j \vee X_k)$  then  $Q_\Phi$  contains  $A(x_i, x_j, x_k)$ .
- If  $C = (X_i \vee X_j \vee \neg X_k)$  then  $Q_\Phi$  contains  $B(x_i, x_j, x_k)$ .
- If  $C = (X_i \vee \neg X_j \vee \neg X_k)$  then  $Q_\Phi$  contains  $C(x_i, x_j, x_k)$ .
- If  $C = (\neg X_i \vee \neg X_j \vee \neg X_k)$  then  $Q_\Phi$  contains  $D(x_i, x_j, x_k)$ .



## Reduction from 3SAT to CQ Evaluation

Given a 3CNF formula  $\Phi$  we construct  $Q_\Phi, \mathbf{D}$  such that:

$\Phi$  is satisfiable iff  $\exists h : Q_\Phi \rightarrow \mathbf{D}$ .

$Q_\Phi$  has one atom for each clause  $C$  in  $\Phi$ :

- If  $C = (X_i \vee X_j \vee X_k)$  then  $Q_\Phi$  contains  $A(x_i, x_j, x_k)$ .
- If  $C = (X_i \vee X_j \vee \neg X_k)$  then  $Q_\Phi$  contains  $B(x_i, x_j, x_k)$ .
- If  $C = (X_i \vee \neg X_j \vee \neg X_k)$  then  $Q_\Phi$  contains  $C(x_i, x_j, x_k)$ .
- If  $C = (\neg X_i \vee \neg X_j \vee \neg X_k)$  then  $Q_\Phi$  contains  $D(x_i, x_j, x_k)$ .

$\mathbf{D}$  has 4 tables with 7 tuples each **which tuple is missing?**

$$A = \begin{array}{|c|c|c|} \hline 0 & 0 & 1 \\ \hline & \vdots & \\ \hline 1 & 1 & 1 \\ \hline \end{array}$$

$$B = \begin{array}{|c|c|c|} \hline 0 & 0 & 0 \\ \hline 0 & 1 & 0 \\ \hline & \vdots & \\ \hline 1 & 1 & 1 \\ \hline \end{array}$$

$C = \dots$

$D = \dots$

## Reduction from 3SAT to CQ Evaluation

Given a 3CNF formula  $\Phi$  we construct  $Q_\Phi, \mathbf{D}$  such that:

$\Phi$  is satisfiable iff  $\exists h : Q_\Phi \rightarrow \mathbf{D}$ .

$Q_\Phi$  has one atom for each clause  $C$  in  $\Phi$ :

- If  $C = (X_i \vee X_j \vee X_k)$  then  $Q_\Phi$  contains  $A(x_i, x_j, x_k)$ .
- If  $C = (X_i \vee X_j \vee \neg X_k)$  then  $Q_\Phi$  contains  $B(x_i, x_j, x_k)$ .
- If  $C = (X_i \vee \neg X_j \vee \neg X_k)$  then  $Q_\Phi$  contains  $C(x_i, x_j, x_k)$ .
- If  $C = (\neg X_i \vee \neg X_j \vee \neg X_k)$  then  $Q_\Phi$  contains  $D(x_i, x_j, x_k)$ .

$\mathbf{D}$  has 4 tables with 7 tuples each **which tuple is missing?**

$$A = \begin{array}{|c|c|c|} \hline 0 & 0 & 1 \\ \hline & \vdots & \\ \hline 1 & 1 & 1 \\ \hline \end{array} \quad B = \begin{array}{|c|c|c|} \hline 0 & 0 & 0 \\ \hline 0 & 1 & 0 \\ \hline & \vdots & \\ \hline 1 & 1 & 1 \\ \hline \end{array} \quad C = \dots \quad D = \dots$$

**In class:**  $\Phi$  is satisfiable iff  $\exists h : Q \rightarrow \mathbf{D}$ .

## Example of 3SAT Reduction

Check if it is satisfiable:

$$\Phi = (X \vee \neg Y \vee Z) \wedge (\neg X \vee Y \vee \neg Z) \wedge (X \vee Y \vee Z) \wedge (\neg X \vee \neg Y \vee Z).$$

## Example of 3SAT Reduction

Check if it is satisfiable:

$$\Phi = (X \vee \neg Y \vee Z) \wedge (\neg X \vee Y \vee \neg Z) \wedge (X \vee Y \vee Z) \wedge (\neg X \vee \neg Y \vee Z).$$

$$Q = B(x, z, y) \wedge C(y, x, z) \wedge A(x, y, z) \wedge C(z, x, y)$$

## Example of 3SAT Reduction

Check if it is satisfiable:

$$\Phi = (X \vee \neg Y \vee Z) \wedge (\neg X \vee Y \vee \neg Z) \wedge (X \vee Y \vee Z) \wedge (\neg X \vee \neg Y \vee Z).$$

$$Q = B(x, z, y) \wedge C(y, x, z) \wedge A(x, y, z) \wedge C(z, x, y)$$

 $A =$ 

0	0	1
0	1	0
0	1	1
1	0	0
1	0	1
1	1	0
1	1	1

 $B =$ 

0	0	0
0	1	0
0	1	1
1	0	0
1	0	1
1	1	0
1	1	1

 $C =$ 

0	0	0
0	0	1
0	1	0
1	0	0
1	0	1
1	1	0
1	1	1

(no need for  $D$ )

## Combined Complexity for FO

Recall that the combined complexity of FO is in PSPACE.

### Theorem

*There exists a database  $\mathbf{D}$  for which the expression complexity of FO queries is PSPACE complete.*

Thus, both expression and combined complexities are PSPACE-complete.

# Combined Complexity for FO

Recall that the combined complexity of FO is in PSPACE.

## Theorem

*There exists a database  $\mathbf{D}$  for which the expression complexity of FO queries is PSPACE complete.*

Thus, both expression and combined complexities are PSPACE-complete.

**Proof:** Reduction from the [Quantified Boolean Formula Satisfiability](#):

$$Q_1 X_1 \ Q_2 X_2 \ \cdots Q_n X_n \ \Phi$$

where  $\Phi$  is 3CNF.

# Combined Complexity for FO

Recall that the combined complexity of FO is in PSPACE.

## Theorem

*There exists a database  $\mathbf{D}$  for which the expression complexity of FO queries is PSPACE complete.*

Thus, both expression and combined complexities are PSPACE-complete.

**Proof:** Reduction from the [Quantified Boolean Formula Satisfiability](#):

$$Q_1 X_1 \ Q_2 X_2 \ \cdots Q_n X_n \ \Phi$$

where  $\Phi$  is 3CNF.

Use the same  $Q_\Phi, \mathbf{D}$  before, but add appropriate quantifiers to  $Q_\Phi$ :

$$Qx_1 \ Qx_2 \ \cdots Q_n X_n \ Q_\Phi(x_1, \dots, x_n)$$



## Example of QBE Reduction

Check if it is true:

$\Phi =$

$$\forall X \exists Y \forall Z ((X \vee \neg Y \vee Z) \wedge (\neg X \vee Y \vee \neg Z) \wedge (X \vee Y \vee Z) \wedge (\neg X \vee \neg Y \vee Z))$$

## Example of QBE Reduction

Check if it is true:

$\Phi =$

$$\forall X \exists Y \forall Z ((X \vee \neg Y \vee Z) \wedge (\neg X \vee Y \vee \neg Z) \wedge (X \vee Y \vee Z) \wedge (\neg X \vee \neg Y \vee Z))$$

$$Q = \forall x \exists y \forall z (B(x, z, y) \wedge C(y, x, z) \wedge A(x, y, z) \wedge C(z, x, y))$$

## Example of QBE Reduction

Check if it is true:

$\Phi =$

$$\forall X \exists Y \forall Z ((X \vee \neg Y \vee Z) \wedge (\neg X \vee Y \vee \neg Z) \wedge (X \vee Y \vee Z) \wedge (\neg X \vee \neg Y \vee Z))$$

$$Q = \forall x \exists y \forall z (B(x, z, y) \wedge C(y, x, z) \wedge A(x, y, z) \wedge C(z, x, y))$$

 $A =$ 

0	0	1
0	1	0
0	1	1
1	0	0
1	0	1
1	1	0
1	1	1

 $B =$ 

0	0	0
0	1	0
0	1	1
1	0	0
1	0	1
1	1	0
1	1	1

 $C =$ 

0	0	0
0	0	1
0	1	0
1	0	0
1	0	1
1	1	0
1	1	1

(no need for  $D$ )

## Discussion: CQ and CSP

The generalized Constraint Satisfaction Problem is:

### Definition (Kolaitis&Vardi)

Given two classes of finite structures  $\mathcal{A}, \mathcal{B}$ , the  $CSP(\mathcal{A}, \mathcal{B})$  problem is:  
Given  $A \in \mathcal{A}, B \in \mathcal{B}$ , is there a homomorphism  $h : A \rightarrow B$ ?

## Discussion: CQ and CSP

The generalized Constraint Satisfaction Problem is:

### Definition (Kolaitis&Vardi)

Given two classes of finite structures  $\mathcal{A}, \mathcal{B}$ , the  $CSP(\mathcal{A}, \mathcal{B})$  problem is:  
Given  $A \in \mathcal{A}, B \in \mathcal{B}$ , is there a homomorphism  $h : A \rightarrow B$ ?

Standard CSP restricts the **right-hand side**,  $CSP(-, B)$ .

What is  $B$  for 3SAT? For 3-colorability? For Hamiltonian cycle?

## Discussion: CQ and CSP

The generalized Constraint Satisfaction Problem is:

### Definition (Kolaitis&Vardi)

Given two classes of finite structures  $\mathcal{A}, \mathcal{B}$ , the  $CSP(\mathcal{A}, \mathcal{B})$  problem is:  
Given  $A \in \mathcal{A}, B \in \mathcal{B}$ , is there a homomorphism  $h : A \rightarrow B$ ?

Standard CSP restricts the **right-hand side**,  $CSP(-, B)$ .

What is  $B$  for 3SAT? For 3-colorability? For Hamiltonian cycle?

Query evaluation restricts the **left-hand side**,  $CSP(Q, -)$

“Query evaluation is **CSP from the other side.**”

# Summary

- Evaluating  $Q(\mathbf{D})$  consists of finding homomorphisms  $h : Q \rightarrow \mathbf{D}$ .
- This problem is in NP, in fact it is the very definition of NP.
- If  $Q$  is fixed, then the problem is in PTIME in  $|\mathbf{D}|$ . [Data complexity](#)
- If  $Q$  is part of the input (i.e. can be huge) then NP-complete.  
[Expression complexity](#)