

Constraint Satisfaction and Database Theory *

Moshe Y. Vardi

Rice University

* Joint work with T. Feder and P. Kolaitis

Constraint Satisfaction Problem (CSP)

Input: (V, D, C) :

- A finite set V of *variables*
- A finite set D of *values*
- A finite set C of *constraints* restricting the values that tuples of variables can take.

Constraint: (t, R)

- t : a tuple of variables over V
- R : a relation of arity $|t|$

Solution: $h : V \rightarrow D$

- $h(t) \in R$: for all $(t, R) \in C$

Question: Does (V, D, C) have a solution? I.e., is there an assignment of values to the variables such that all constraints are satisfied?

3-Colorability

3-COLOR: Given an undirected graph $A = (V, E)$, is it 3-colorable?

- The variables are the nodes in V .
- The values are the elements in $\{\mathbf{R}, \mathbf{G}, \mathbf{B}\}$.
- The constraints are $\{(\langle u, v \rangle, \rho) : (u, v) \in E\}$, where $\rho = \{(R, G), (R, B), (G, R), (G, B), (B, R), (B, G)\}$.

Introduction to Database Theory

Basic Concepts:

- *Relation Scheme*: a set of attributes
- *Tuple*: mapping from relation scheme to data values
- *Tuple Projection*: if t is a tuple on P , and $Q \subseteq P$, then $t[Q]$ is the restriction of t to Q .
- *Relation*: a set of tuples over a relation scheme
- *Relational Projection*: if R is a relation on P , and $Q \subseteq P$, then $R[Q]$ is the relation $\{t[Q] : t \in R\}$.
- *Join*: Let R_i be a relation over relation scheme S_i . Then $\bowtie_i R_i$ is a relation over the relation scheme $\cup_i S_i$ defined by $\bowtie_i R_i = \{t : t[S_i] \in R_i\}$.

Database Perspective of CSP

Given: $(V, D, \{C_1, \dots, C_m\})$, where $C_i = (t_i, R_i)$.

Assume (wlog): Each t_i consists of distinct elements.

Database Perspective:

- V : attributes
- D : values
- (t_i, R_i) : relation R_i over relation scheme t_i

Fact: (Bibel, Gyssens, Jeavons, Cohen)

$(V, D, \{C_1, \dots, C_m\})$ has a solution iff $\bowtie_1^m R_i$ is nonempty.

Homomorphisms

Homomorphism: Let $\mathbf{A} = (A, R_1^{\mathbf{A}}, \dots, R_m^{\mathbf{A}})$ and $\mathbf{B} = (B, R_1^{\mathbf{B}}, \dots, R_m^{\mathbf{B}})$ be two relational structures.

$h : A \rightarrow B$ is a *homomorphism* from \mathbf{A} to \mathbf{B} if for every $i \leq m$ and every tuple $(a_1, \dots, a_n) \in A^n$,

$$R_i^{\mathbf{A}}(a_1, \dots, a_n) \implies R_i^{\mathbf{B}}(h(a_1), \dots, h(a_n)).$$

The Homomorphism Problem: Given relational structures \mathbf{A} and \mathbf{B} , is there a homomorphism $h : \mathbf{A} \rightarrow \mathbf{B}$?

Example: An undirected graph $\mathbf{A} = (V, E)$ is 3-colorable



there is a homomorphism $h : \mathbf{A} \rightarrow K_3$, where K_3 is the *3-clique*.

Homomorphism Problems

Examples:

- k -Clique: $K_k \xrightarrow{h} (V, E)?$
- Hamiltonian Cycle: $(V, C_{|V|}, \neq) \xrightarrow{h} (V, E, \neq)?$
- Subgraph Isomorphism: $(V, E, \overline{E}) \xrightarrow{h} (V', E', \overline{E'})?$
- s - t Connectivity: $(V, E, \{\langle s, t \rangle\}) \xrightarrow{h} (\{0, 1\}, =, \neq)?$

Fact: (Levin, 1973)

The homomorphism problem is NP-complete.

CSP vs. Homomorphisms

From CSP to Homomorphism:

Given: $(V, D, \{C_1, \dots, C_m\})$, where $C_i = (t_i, R_i)$.

Define **A**, **B**:

- $\mathbf{A} = (V, \{t_1\}, \dots, \{t_m\})$
- $\mathbf{B} = (D, R_1, \dots, R_m)$

Fact: (V, D, C) has a solution iff there is homomorphism from **A** to **B**.

CSP vs. Homomorphisms

From Homomorphism to CSP:

Given: $\mathbf{A} = (A, R_1^{\mathbf{A}}, \dots, R_m^{\mathbf{A}})$, $\mathbf{B} = (B, R_1^{\mathbf{B}}, \dots, R_m^{\mathbf{B}})$.

Define (V, D, C) :

- $V = A$: elements of \mathbf{A} are variables.
- $D = B$: elements of \mathbf{B} are values.
- $C = \{(t, R_i^{\mathbf{B}}) : t \in R_i^{\mathbf{A}}\}$: constraints derived from \mathbf{A}, \mathbf{B} .

Fact: There is homomorphism from \mathbf{A} to \mathbf{B} iff (V, D, C) has a solution.

Conclusion: CSP=Homomorphism Problem

- Feder&V., 1993
- Garey&Johnson, 1979: Homomorphism in, CSP not.

Conjunctive Queries

Conjunctive Query: First-order logic without \forall, \exists, \neg ;
written as a rule

$$Q(X_1, \dots, X_n) : - R(X_3, Y_2, X_4), \dots, S(X_2, Y_3)$$

Significance: most common SQL queries (*Select-Project-Join*)

Example:

$$GrandParent(X, Y) : - Parent(X, Z), Parent(Z, Y)$$

Conjunctive Query Containment (CQC)

Definition: $Q_1 \sqsubseteq Q_2$ iff $Q_1(B) \subseteq Q_2(B)$ for every database B .

Fundamental problem in database theory:

- Optimization of SPJ Queries
- Information Integration
- Data Warehousing
- ...

CQC vs. Homomorphism Problem

Canonical Database B^Q :

- Each variable in Q is a distinct element
- Each subgoal $R(X_3, Y_2, X_4)$ of Q gives rise to a tuple $R(X_3, Y_2, X_4)$ in B^Q
- Each *distinguished variable* X_i gives rise to a tuple $P_i(X_i)$ in B^Q

Fact: (Chandra&Merlin, 1977)

For conjunctive queries Q_1 and Q_2 , TFAE:

- $Q_1 \sqsubseteq Q_2$
- There is a homomorphism $h : B^{Q_2} \rightarrow B^{Q_1}$.
- $Q_2(B^{Q_1})$ is nonempty

Homomorphism Problem vs. CQC

Given: $\mathbf{A} = (A, R_1^A, \dots, R_m^A)$, $\mathbf{B} = (B, R_1^B, \dots, R_m^B)$.

Define conjunctive query Q_A over variables in A :

$$Q_A : - \bigwedge_{i=1}^m \bigwedge_{t \in R_i^A} R_i(t)$$

Define conjunctive query Q_B over variables in B :

$$Q_B : - \bigwedge_{i=1}^m \bigwedge_{t \in R_i^B} R_i(t)$$

Fact: For relational structures \mathbf{A} and \mathbf{B} , TFAE:

- There is a homomorphism $h : \mathbf{A} \rightarrow \mathbf{B}$.
- $Q_B \sqsubseteq Q_A$
- $Q_A(\mathbf{B})$ is nonempty

“Same Difference”

Conclusion: (Kolaitis&V., 1998)

CSP, CQ containment, and CQ evaluation amount to the same problem!

Potential Promise: Synergy between AI and DB.

Uniform CSP vs. Non-Uniform CSP

Uniform CSP:

$$\{(\mathbf{A}, \mathbf{B}) : \exists \text{ homomorphism } h : \mathbf{A} \rightarrow \mathbf{B}\}$$

Complexity of Uniform CSP: NP-complete

Non-uniform CSP: Fix a structure \mathbf{B}

$$\text{CSP}(\mathbf{B}) = \{\mathbf{A} : \exists \text{ homomorphism } h : \mathbf{A} \rightarrow \mathbf{B}\}$$

Complexity of Non-Uniform CSP: Depends on \mathbf{B}

- $\text{CSP}(K_2)$ is in PTIME (2-COLORABILITY)
- $\text{CSP}(K_3)$ is NP-complete (3-COLORABILITY)

Complexity of Query Evaluation

Measuring Complexity:

- *Combined Complexity:*

$$\langle \{Q, \mathbf{B}\} : Q(\mathbf{B}) \text{ is nonempty} \rangle$$

- *Expression Complexity:* Fix \mathbf{B}

$$\{Q : Q(\mathbf{B}) \text{ is nonempty}\}$$

- *Data Complexity:* Fix Q

$$\{\mathbf{B} : Q(\mathbf{B}) \text{ is nonempty}\}$$

Database Perspective

Uniform CSP:

$$\{(\mathbf{A}, \mathbf{B}) : Q_{\mathbf{A}}(\mathbf{B}) \text{ is nonempty}\}$$

Complexity of Uniform CSP: combined complexity

Non-uniform CSP: Fix a structure \mathbf{B}

$$\text{CSP}(\mathbf{B}) = \{\mathbf{A} : Q_{\mathbf{A}}(\mathbf{B}) \text{ is nonempty}\}$$

Complexity of Non-Uniform CSP: expression complexity

Complexity of CQC

In Conjunctive Query Containment ($Q_1 \sqsubseteq Q_2?$), we are typically interested in *uniform* results.

Research Program:

Identify the tractable cases of uniform CQC.

Results: Tractable cases of $Q_1 \sqsubseteq Q_2$

- Saraiya, 1991:

Every predicate occurs at most *twice* in Q_1

- Chekuri–Rajaraman, 1997:

Q_2 has *querywidth* at most k , where k is a fixed positive integer.

- Gottlob–Leone–Scarcello, 1999:

Q_2 has *hypertreewidth* at most k , where k is a fixed positive integer.

Complexity of Non-Uniform CSP

Research Program:

Identify the tractable cases of non-uniform CSP (Schaefer, 1978, ..., Jeavons et al)

Dichotomy Conjecture: (Feder–V., 1993)

For every structure \mathbf{B} ,

- either $\text{CSP}(\mathbf{B})$ is in PTIME
- or $\text{CSP}(\mathbf{B})$ is NP-complete.

Recall: $P \neq NP \Rightarrow NP - NPC - P \neq \emptyset$ (Ladner, 1975)

“Evidence”: (Hell–Nešetřil, 1990)

Let \mathbf{B} be an *undirected* graph.

- \mathbf{B} bipartite \implies $\text{CSP}(\mathbf{B})$ is in PTIME
- \mathbf{B} non-bipartite \implies $\text{CSP}(\mathbf{B})$ is NP-complete

Complexity of Non-Uniform CSP

More “Evidence”: (Schaefer, 1978)

Let \mathbf{B} have a *Boolean* domain, then

- either $\text{CSP}(\mathbf{B})$ is in PTIME
- or $\text{CSP}(\mathbf{B})$ is NP-complete.

Classification Question:

For a given structure \mathbf{B} ,

- when is $\text{CSP}(\mathbf{B})$ in PTIME?
- when is $\text{CSP}(\mathbf{B})$ NP-complete?

Classification Conjecture: (Feder–V., 1993)

Two explanations for tractability of $\text{CSP}(\mathbf{B})$

- *combinatorial* (Datalog)
- *algebraic* (group-theoretic)

Datalog and Non-Uniform CSP

Example: NON 2-COLORABILITY

$$O(X, Y) : - E(X, Y)$$

$$O(X, Y) : - O(X, Z), E(Z, W), E(W, Y)$$

$$Q : - O(X, X)$$

Recall: Datalog \subseteq PTIME

Define: $\overline{\text{CSP}(\mathbf{B})} = \{\mathbf{A} : \mathbf{A} \notin \text{CSP}(\mathbf{B})\}$.

Datalog vs. Non-Uniform CSP: Explanation for many tractability results

- $\overline{\text{CSP}(\mathbf{B})}$ is expressible in Datalog

Note: $\overline{\text{CSP}(\mathbf{B})}$ is positively monotone.

k -Datalog

Definition:

- k -Datalog: Datalog with at most k variables per rule (NON 2-COLORABILITY is in 4-Datalog)
- $\exists\text{IL}^k$: k -variable existential positive infinitary logic
 - variables: x_1, \dots, x_k
 - no universal quantifiers
 - no negations
 - infinitary conjunctions and disjunctions

Facts: Fix $k \geq 1$

- k -Datalog $\subset \exists\text{IL}^k$
- $\exists\text{IL}^k$ can be characterized in terms of *existential k -pebble games* between the *Spoiler* and the *Duplicator*.
- There is a PTIME algorithm to decide whether the *Spoiler* or the *Duplicator* wins the existential k -pebble game.

Existential k -Pebble Games

A, B: structures

- **Spoiler:** *places on or removes* a pebble from an element of **A**.
- **Duplicator:** tries to duplicate move on **B**.

A: a_1, a_2, \dots, a_l $l \leq k$

B: b_1, b_2, \dots, b_l

- *Spoiler wins:* $h(a_i) = b_i, 1 \leq i \leq l$ is **not** a homomorphism.
- *Duplicator wins:* otherwise.

Fact: (Kolaitis&V., 1995)

B satisfies the same $\exists\text{IL}^k$ sentences as **A** iff the Duplicator wins the existential k -pebble game on **A, B**.

k -Datalog and CSP

Theorem: (Kolaitis&V., 1998): TFAE for $k \geq 1$ and a structure \mathbf{B} :

- $\overline{\text{CSP}(\mathbf{B})}$ is expressible in k -Datalog
- $\overline{\text{CSP}(\mathbf{B})}$ is expressible in $\exists\text{IL}^k$
- $\text{CSP}(\mathbf{B}) = \{\mathbf{A} : \text{Duplicator wins the existential } k\text{-pebble game on } \mathbf{A} \text{ and } \mathbf{B}\}.$

Intuition: $\overline{\text{CSP}(\mathbf{B})} \in k\text{-Datalog}$ implies that existence of homomorphism is equivalent to the Duplicator winning the existential k -pebble game.

Winning Configurations

Let $\mathbf{A} = (A, R_1^{\mathbf{A}}, \dots, R_m^{\mathbf{A}})$ and $\mathbf{B} = (B, R_1^{\mathbf{B}}, \dots, R_m^{\mathbf{B}})$ be two relational structures. Let $\mathbf{a} \in A^k$ and $\mathbf{b} \in B^k$.

Definition:

- (\mathbf{a}, \mathbf{b}) is a *winning configuration* for the Duplicator if the Duplicator wins the existential k -pebble game starting from the position (\mathbf{a}, \mathbf{b}) .
- $W^k(\mathbf{A}, \mathbf{B})$ is the set of winning configurations.

Theorem: (Kolaitis&V., 1998)

- $W^k(\mathbf{A}, \mathbf{B})$ is definable in least-fixpoint logic over the combined structure $\mathbf{A} + \mathbf{B}$.
- For a fixed structure \mathbf{B} , the complement of $W^k(\mathbf{A}, \mathbf{B})$ is definable in k -Datalog over the structure \mathbf{A} .

Corollary: For a fixed structure \mathbf{B} , there is a k -Datalog program $\rho_{\mathbf{B}}$ such that $\rho_{\mathbf{B}}(\mathbf{A})$ is nonempty iff the Spoiler wins the existential k -pebble game on \mathbf{A}, \mathbf{B} .

k -Datalog and CSP

$\rho_{\mathbf{B}}$:

- If $\rho_{\mathbf{B}}(\mathbf{A})$ is nonempty, then $\mathbf{A} \notin \text{CSP}(\mathbf{B})$.
- If $\overline{\text{CSP}(\mathbf{B})}$ is definable in k -Datalog, then it is definable by $\rho_{\mathbf{B}}$.
- *Open question:* Decide for a given \mathbf{B} whether $\overline{\text{CSP}(\mathbf{B})}$ is definable by $\rho_{\mathbf{B}}$.

Database Perspective of k -Datalog: Fix \mathbf{B}

$$\text{CSP}(\mathbf{B}) = \{\mathbf{A} : \rho_{\mathbf{B}}(\mathbf{A}) \text{ is empty}\}$$

Complexity of Non-Uniform CSP: data complexity of Datalog

Earlier Database Perspective

Non-uniform CSP: Fix a structure \mathbf{B}

$$\text{CSP}(\mathbf{B}) = \{\mathbf{A} : Q_{\mathbf{A}}(\mathbf{B}) \text{ is nonempty}\}$$

Complexity of Non-Uniform CSP: expression
complexity of conjunctive queries

Winning Strategies

\mathbf{A}, \mathbf{B} : structures

Definition: A *winning strategy* for the Duplicator in the existential k -pebble game on \mathbf{A}, \mathbf{B} is a nonempty family \mathcal{F} of k -partial homomorphisms from \mathbf{A} to \mathbf{B} that has the *k -forth property*:

- for every $f \in \mathcal{F}$ with $|f| < k$ and every $a \in A$ on which f is undefined, there is a $g \in \mathcal{F}$ that extends f and is defined on a .

Intuition:

- Think of f as a position in the game.
- If Spoiler places pebble on a , Duplicators responds with a pebble on $g(a)$.

CSP and Consistency

Key Idea in AI Research: *consistency*

- extensibility of partial solutions
- backtrack-free search

Definition: A homomorphism instance \mathbf{A}, \mathbf{B} is

- *i -consistent*: For every consistent assignment α to $i - 1$ variables v_1, \dots, v_{i-1} and a variable v_i , there is a consistent extension of α to v_i .
- *Strongly k -consistent*: i -consistency holds for all $i \leq k$.

Game-Theoretic Formulation: (Kolaitis&V., 2000)

\mathbf{A}, \mathbf{B} is strongly k -consistent iff the family of all k -partial homomorphism from \mathbf{A} to \mathbf{B} is a winning strategy for the Duplicator in the existential k -pebble game on \mathbf{A}, \mathbf{B} .

From Local to Global Consistency

Condition for Tractability: (Dechter, 1992)

Establishing strong k -consistency implies existence of solution.

- Start with \mathbf{A}, \mathbf{B} .
- Generate a strongly k -consistent instance \mathbf{A}', \mathbf{B}' over the same domains.
- If h is a k -partial homomorphism from \mathbf{A}' to \mathbf{B}' , then h is a k -partial homomorphism from \mathbf{A} to \mathbf{B} .
- A function h is a homomorphism from \mathbf{A} to \mathbf{B} iff it is a homomorphism from \mathbf{A}' to \mathbf{B}' .

Intuition: \mathbf{A}', \mathbf{B}' is obtained by inferring from \mathbf{A}, \mathbf{B} all constraints on k -variables – *constraint propagation*.

Consistency and k -Datalog

Strong k -Consistency vs. Games: (K.&V., 2000)

The strongly k -consistent instance \mathbf{A}', \mathbf{B}' can be obtained by “re-formatting” $W^k(\mathbf{A}, \mathbf{B})$

Intuition: Computing winning strategies can be viewed as constraint propagation.

Corollary: Let \mathbf{B} a relational structure. $\overline{\text{CSP}(\mathbf{B})}$ is expressible in k -Datalog iff for every structure \mathbf{A} , establishing strong k -consistency for \mathbf{A}, \mathbf{B} implies that there is a homomorphism from \mathbf{A} to \mathbf{B} .

Constraint Propagation

A database perspective on constraint propagation:

- *Projection Rule*: From (t, R) infer $(t, \pi_{t'}(R))$, where $t' \subseteq t$.
- *Join Rule*: From (t_1, R_1) and (t_2, R_2) , infer $(t_1 \cup t_2, R_1 \bowtie R_2)$

Example: *Resolution* – From $\alpha = \theta \vee p$ and $\beta = \psi \vee \neg p$ infer $\theta \vee \psi$

- Sound and complete for refutation of clausal formulas

Soundness Proof:

- Consider $models(\varphi)$ as a relation (over a Boolean domain) with relation scheme $vars(\varphi)$.
- $(models(\alpha) \bowtie models(\beta)) [vars(\theta) \cup vars(\psi)] = models(\theta \vee \psi)$

Bounded Treewidth

Definition: A *tree decomposition* of a structure $\mathbf{A} = (A, R_1, \dots, R_m)$ is a labeled tree T such that

- Each label is a non-empty subset of A ;
- For every R_i and every $(a_1, \dots, a_n) \in R_i$, there is a node whose label contains $\{a_1, \dots, a_n\}$.
- For every $a \in A$, the nodes whose label contain a form a subtree.

The *treewidth* $\text{tw}(\mathbf{A})$ of \mathbf{A} is defined by

$$\text{tw}(\mathbf{A}) = \min_T \{ \max \{ \text{label size in } T \} \} - 1$$

Note: Generalizes the *treewidth* of a *graph*.

Bounded Treewidth and CSP

$$C_k = \{\mathbf{A} : \text{tw}(\mathbf{A}) \leq k\}$$

Theorem: (Freuder, 1990)

$\text{CSP}(C_k, \text{All})$ is in PTIME.

Note:

- Complexity is exponential in k .
- Determining treewidth of \mathbf{B} is NP-hard.
- Checking if treewidth is k is in linear time.

Database Perspective of Bounded Treewidth

Tractable cases: CSP(A, B)

- Yannakakis, 1981:

A is *acyclic*

- Freuder, 1990:

A has *treewidth* at most k

- Chekuri–Rajaraman, 1997:

A has *querywidth* at most k

- Gottlob–Leone–Scarcello, 1999:

A has *hypertreewidth* at most k

Acyclicity is the crux: (Flum–Frick–Grohe, 2000)

The above can all be viewed as instances of Yannakakis's algorithm.

Complexity of Query Evaluation

Expression Complexity: Fix \mathbf{B}

$$\{Q : Q(\mathbf{B}) \text{ is nonempty}\}$$

Data Complexity: Fix Q

$$\{\mathbf{B} : Q(\mathbf{B}) \text{ is nonempty}\}$$

Exponential Gap: (V., 1982)

- Data complexity of FO: *LOGSPACE*
- Expression complexity of FO: *PSPACE-complete*

Mystery: practical query evaluation

Variable-Confined Queries

Definition: FO^k is first-order logic with at most k variables.

In Practice: (V., 1995)

- Queries often can be rewritten to use a small number of variables.
- Variable-confined queries have lower expression complexity.
- E.g.: expression complexity of FO^k is *PTIME-complete*

Recall: $\mathbf{A} \mapsto Q_{\mathbf{A}}$

- $h : \mathbf{A} \rightarrow \mathbf{B}$ iff $Q_{\mathbf{A}}(\mathbf{B})$ is nonempty

Theorem: (Kolaitis&V., 1998)

$Q_{\mathbf{A}}$ is equivalent to a CQ of treewidth k iff $Q_{\mathbf{A}}$ is expressible in existential, positive FO with $k + 1$ variables.

AI vs. Databases

Fundamentally different approaches:

- *AI*: search one answer
- *Databases*: generate all answers

Fundamental Question: Is query evaluation effective for solving CSP?

Recall: $\mathbf{A} \mapsto Q_{\mathbf{A}}$

- $h : \mathbf{A} \rightarrow \mathbf{B}$ iff $Q_{\mathbf{A}}(\mathbf{B})$ is nonempty

Experiment: (with Alfonso San Miguel Aguirre)

- Express 3-SAT as a CSP.
- Translate to conjunctive query evaluation.
- Run on DB2.

Outcome: Terrible performance!!!

- DB2 optimizer ill suited.

In Conclusion

CSP is a database problem!!!

- CSP vs. conjunctive queries
- CSP vs. Datalog
- CSP vs. variable-confined queries

There is a real potential for AI-DB cross fertilization.