# DATABASE THEORY

## Lecture 14: Datalog Evaluation

**Markus Krötzsch**

**Knowledge-Based Systems**

TU Dresden, 16 June 2025

# Review: Datalog

A rule-based recursive query language

---

father(alice, bob)

mother(alice, carla)

$\text{Parent}(x, y) \leftarrow \text{father}(x, y)$

$\text{Parent}(x, y) \leftarrow \text{mother}(x, y)$

SameGeneration(x, x)

$\text{SameGeneration}(x, y) \leftarrow \text{Parent}(x, v) \land \text{Parent}(y, w) \land \text{SameGeneration}(v, w)$

---

- Datalog is more complex than FO query answering
- Datalog is more expressive than FO query answering
- Semipositive Datalog with a successor ordering captures P
- Datalog containment is undecidable

Remaining question: How can Datalog query answering be implemented?

# Implementing Datalog

FO queries (and thus also CQs and UCQs) are supported by almost all DBMS
⤳ many specific implementation and optimisation techniques

How can Datalog queries be answered in practice?
⤳ techniques for dealing with recursion in DBMS query answering

There are two major paradigms for answering recursive queries:

- Bottom-up: derive conclusions by applying rules to given facts
- Top-down: search for proofs to infer results given query

# Computing Datalog Query Answers Bottom-Up

We already saw a way to compute Datalog answers bottom-up:
the step-wise computation of the consequence operator $T_P$

Bottom-up computation is known under many names:

- Forward-chaining since rules are "chained" from premise to conclusion
  (common in logic programming)

- Materialisation since inferred facts are stored ("materialised")
  (common in databases)

- Saturation since the input database is "saturated" with inferences
  (common in theorem proving)

- Deductive closure since we "close" the input under entailments
  (common in formal logic)

# Naive Evaluation of Datalog Queries

A direct approach for computing $T_P^\infty$

| | |
|---|---|
| 01 | $T_P^0 := \emptyset$ |
| 02 | $i := 0$ |
| 03 | **repeat** : |
| 04 | $\quad T_P^{i+1} := \emptyset$ |
| 05 | $\quad$ **for** $H \leftarrow B_1 \wedge \ldots \wedge B_\ell \in P$ : |
| 06 | $\quad\quad$ **for** $\theta \in B_1 \wedge \ldots \wedge B_\ell(T_P^i)$ : |
| 07 | $\quad\quad\quad T_P^{i+1} := T_P^{i+1} \cup \{H\theta\}$ |
| 08 | $\quad i := i + 1$ |
| 09 | **until** $T_P^{i-1} = T_P^i$ |
| 10 | **return** $T_P^i$ |

Notation for line 06/07:

- a substitution $\theta$ is a mapping from variables to database elements

- for a formula $F$, we write $F\theta$ for the formula obtained by replacing each free variable $x$ in $F$ by $\theta(x)$

- for a CQ $Q$ and database $\mathcal{I}$, we write $\theta \in Q(\mathcal{I})$ if $\mathcal{I} \models Q\theta$

# What's Wrong with Naive Evaluation?

An example Datalog program:

$$e(1, 2) \quad e(2, 3) \quad e(3, 4) \quad e(4, 5)$$

$$(R1) \qquad T(x, y) \leftarrow e(x, y)$$

$$(R2) \qquad T(x, z) \leftarrow T(x, y) \land T(y, z)$$

How many body matches do we need to iterate over?

$T_P^0 = \emptyset$          initialisation

$T_P^1 = \{T(1, 2), T(2, 3), T(3, 4), T(4, 5)\}$     4 matches for $(R1)$

$T_P^2 = T_P^1 \cup \{T(1, 3), T(2, 4), T(3, 5)\}$      $4 \times (R1) + 3 \times (R2)$

$T_P^3 = T_P^2 \cup \{T(1, 4), T(2, 5), T(1, 5)\}$      $4 \times (R1) + 8 \times (R2)$

$T_P^4 = T_P^3 = T_P^\infty$               $4 \times (R1) + 10 \times (R2)$

In total, we considered 37 matches to derive 11 facts

# Less Naive Evaluation Strategies

Does it really matter how often we consider a rule match?
After all, each fact is added only once ...

In practice, finding applicable rules takes significant time, even if the conclusion does not need to be added – iteration takes time!
⤳ huge potential for optimisation

**Observation:**
we derive the same conclusions over and over again in each step

**Idea:** apply rules only to newly derived facts
⤳ semi-naive evaluation

## Semi-Naive Evaluation

The computation yields sets $T_P^0 \subseteq T_P^1 \subseteq T_P^2 \subseteq \ldots \subseteq T_P^\infty$

- For an IDB predicate R, let $R^i$ be the "predicate" that contains exactly the R-facts in $T_P^i$
- For $i \leq 1$, let $\Delta_R^i$ be the collection of facts $R^i \setminus R^{i-1}$

We can restrict rules to use only some computations.

**Some options for the computation in step $i + 1$:**

$$\mathsf{T}(x,z) \leftarrow \mathsf{T}^i(x,y) \wedge \mathsf{T}^i(y,z) \qquad \text{same as original rule}$$
$$\mathsf{T}(x,z) \leftarrow \Delta_\mathsf{T}^i(x,y) \wedge \Delta_\mathsf{T}^i(y,z) \qquad \text{restrict to new facts}$$
$$\mathsf{T}(x,z) \leftarrow \Delta_\mathsf{T}^i(x,y) \wedge \mathsf{T}^i(y,z) \qquad \text{partially restrict to new facts}$$
$$\mathsf{T}(x,z) \leftarrow \mathsf{T}^i(x,y) \wedge \Delta_\mathsf{T}^i(y,z) \qquad \text{partially restrict to new facts}$$

What to choose?

# Semi-Naive Evaluation (2)

Inferences that involve new and old facts are necessary:

$$e(1, 2) \quad e(2, 3) \quad e(3, 4) \quad e(4, 5)$$

$$(R1) \qquad T(x, y) \leftarrow e(x, y)$$

$$(R2) \qquad T(x, z) \leftarrow T(x, y) \wedge T(y, z)$$

$$T_P^0 = \emptyset$$

$$\Delta_T^1 = \{T(1, 2), T(2, 3), T(3, 4), T(4, 5)\} \qquad T_P^1 = \Delta_T^1$$

$$\Delta_T^2 = \{T(1, 3), T(2, 4), T(3, 5)\} \qquad T_P^2 = T_P^1 \cup \Delta_T^2$$

$$\Delta_T^3 = \{T(1, 4), T(2, 5), T(1, 5)\} \qquad T_P^3 = T_P^2 \cup \Delta_T^3$$

$$\Delta_T^4 = \emptyset \qquad T_P^4 = T_P^3 = T_P^\infty$$

To derive $T(1, 4)$ in $\Delta_T^3$, we need to combine
$T(1, 3) \in \Delta_T^2$ with $T(3, 4) \in \Delta_T^1$ or $T(1, 2) \in \Delta_T^1$ with $T(2, 4) \in \Delta_T^2$
$\leadsto$ rule $T(x, z) \leftarrow \Delta_T^i(x, y) \wedge \Delta_T^i(y, z)$ is not enough

## Semi-Naive Evaluation (3)

**Correct approach:** consider only rule application that use at least one newly derived IDB atom

For example program:

$$e(1, 2) \quad e(2, 3) \quad e(3, 4) \quad e(4, 5)$$

$$(R1) \qquad T(x, y) \leftarrow e(x, y)$$

$$(R2.1) \qquad T(x, z) \leftarrow \Delta_T^i(x, y) \wedge T^i(y, z)$$

$$(R2.2) \qquad T(x, z) \leftarrow T^i(x, y) \wedge \Delta_T^i(y, z)$$

There is still redundancy here: the matches for $T(x, z) \leftarrow \Delta_T^i(x, y) \wedge \Delta_T^i(y, z)$ are covered by both $(R2.1)$ and $(R2.2)$

$\rightsquigarrow$ replace $(R2.2)$ by the following rule:

$$(R2.2') \qquad T(x, z) \leftarrow T^{i-1}(x, y) \wedge \Delta_T^i(y, z)$$

EDB atoms do not change, so their $\Delta$ would be $\emptyset$

$\rightsquigarrow$ ignore such rules after the first iteration

$$e(1,2) \quad e(2,3) \quad e(3,4) \quad e(4,5)$$

$$(R1) \qquad T(x,y) \leftarrow e(x,y)$$

$$(R2.1) \qquad T(x,z) \leftarrow \Delta_T^i(x,y) \wedge T^i(y,z)$$

$$(R2.2') \qquad T(x,z) \leftarrow T^{i-1}(x,y) \wedge \Delta_T^i(y,z)$$

How many body matches do we need to iterate over?

$$T_P^0 = \emptyset \qquad\qquad\qquad\qquad\qquad\qquad \text{initialisation}$$

$$T_P^1 = \{T(1,2), T(2,3), T(3,4), T(4,5)\} \quad 4 \times (R1)$$

$$T_P^2 = T_P^1 \cup \{T(1,3), T(2,4), T(3,5)\} \qquad 3 \times (R2.1)$$

$$T_P^3 = T_P^2 \cup \{T(1,4), T(2,5), T(1,5)\} \qquad 3 \times (R2.1), 2 \times (R2.2')$$

$$T_P^4 = T_P^3 = T_P^\infty \qquad\qquad\qquad\qquad 1 \times (R2.1), 1 \times (R2.2')$$

In total, we considered 14 matches to derive 11 facts

# Semi-Naive Evaluation: Full Definition

In general, a rule of the form

$$H(\vec{x}) \leftarrow e_1(\vec{y}_1) \wedge \ldots \wedge e_n(\vec{y}_n) \wedge I_1(\vec{z}_1) \wedge I_2(\vec{z}_2) \wedge \ldots \wedge I_m(\vec{z}_m)$$

is transformed into $m$ rules

$$H(\vec{x}) \leftarrow e_1(\vec{y}_1) \wedge \ldots \wedge e_n(\vec{y}_n) \wedge \Delta_{I_1}^i(\vec{z}_1) \wedge I_2^i(\vec{z}_2) \wedge \ldots \wedge I_m^i(\vec{z}_m)$$
$$H(\vec{x}) \leftarrow e_1(\vec{y}_1) \wedge \ldots \wedge e_n(\vec{y}_n) \wedge I_1^{i-1}(\vec{z}_1) \wedge \Delta_{I_2}^i(\vec{z}_2) \wedge \ldots \wedge I_m^i(\vec{z}_m)$$
$$\ldots$$
$$H(\vec{x}) \leftarrow e_1(\vec{y}_1) \wedge \ldots \wedge e_n(\vec{y}_n) \wedge I_1^{i-1}(\vec{z}_1) \wedge I_2^{i-1}(\vec{z}_2) \wedge \ldots \wedge \Delta_{I_m}^i(\vec{z}_m)$$

**Advantages and disadvantages:**

- Huge improvement over naive evaluation
- Some redundant computations remain (see example)
- Some overhead for implementation (store level of entailments)

# Goal-Directed Datalog Evaluation

# Top-Down Evaluation

**Idea:** we may not need to compute all derivations to answer a particular query

---

**Example 14.1:**

$$e(1,2) \quad e(2,3) \quad e(3,4) \quad e(4,5)$$
$$(R1) \quad T(x,y) \leftarrow e(x,y)$$
$$(R2) \quad T(x,z) \leftarrow T(x,y) \wedge T(y,z)$$
$$Query(z) \leftarrow T(2,z)$$

The answers to Query are the T-successors of $2$.

However, bottom-up computation would also produce facts like $T(1,4)$, which are neither directly nor indirectly relevant for computing the query result.

---

# Assumption

**Assumption:** For all techniques presented in this lecture, we assume that the given Datalog program is safe.

- This is without loss of generality (as shown in exercise).
- One can avoid this by adding more cases to algorithms.

# Query-Subquery (QSQ)

QSQ is a technique for organising top-down Datalog query evaluation

**Main principles:**

- Apply backward chaining/resolution: start with query, find rules that can derive query, evaluate body atoms of those rules (subqueries) recursively
- Evaluate intermediate results "set-at-a-time" (using relational algebra on tables)
- Evaluate queries in a "data-driven" way, where operations are applied only to newly computed intermediate results (similar to idea in semi-naive evaluation)
- "Push" variable bindings (constants) from heads (queries) into bodies (subqueries)
- "Pass" variable bindings (constants) "sideways" from one body atom to the next

Details can be realised in several ways.

# Adornments

To guide evaluation, we distinguish free and bound parameters in a predicate.

> **Example 14.2:** If we want to derive atom $T(2, z)$ from the rule $T(x, z) \leftarrow T(x, y) \wedge T(y, z)$, then $x$ will be bound to $2$, while $z$ is free.

We use adornments to denote the free/bound parameters in predicates.

> **Example 14.3:**
>
> $$T^{bf}(x, z) \leftarrow T^{bf}(x, y) \wedge T^{bf}(y, z)$$
>
> - since $x$ is bound in the head, it is also bound in the first atom
> - any match for the first atom binds $y$, so $y$ is bound when evaluating the second atom (in left-to-right evaluation)

# Adornments: Examples

The adornment of the head of a rule determines the adornments of the body atoms:

$$R^{bbb}(x, y, z) \leftarrow R^{bbf}(x, y, v) \wedge R^{bbb}(x, v, z)$$
$$R^{fbf}(x, y, z) \leftarrow R^{fbf}(x, y, v) \wedge R^{bbf}(x, v, z)$$

The order of body predicates affects the adornment:

$$S^{fff}(x, y, z) \leftarrow T^{ff}(x, v) \wedge T^{ff}(y, w) \wedge R^{bbf}(v, w, z)$$
$$S^{fff}(x, y, z) \leftarrow R^{fff}(v, w, z) \wedge T^{fb}(x, v) \wedge T^{fb}(y, w)$$

$\rightsquigarrow$ For optimisation, some orders might be better than others

# Auxiliary Relations for QSQ

To control evaluation, we store intermediate results in auxiliary relations.

When we "call" a rule with a head where some variables are bound, we need to provide the bindings as input
$\leadsto$ for adorned relation $R^{\alpha}$, we use an auxiliary relation $\text{input}_R^{\alpha}$
$\leadsto$ arity of $\text{input}_R^{\alpha}$ = number of $b$ in $\alpha$

The result of calling a rule should be the "completed" input, with values for the unbound variables added
$\leadsto$ for adorned relation $R^{\alpha}$, we use an auxiliary relation $\text{output}_R^{\alpha}$
$\leadsto$ arity of $\text{output}_R^{\alpha}$ = arity of R (= length of $\alpha$)

## Auxiliary Relations for QSQ (2)

When evaluating body atoms from left to right, we use supplementary relations $\text{sup}_i$

$\leadsto$ bindings required to evaluate rest of rule after the $i$th body atom

$\leadsto$ the first set of bindings $\text{sup}_0$ comes from $\text{input}_R^\alpha$

$\leadsto$ the last set of bindings $\text{sup}_n$ go to $\text{output}_R^\alpha$

---

**Example 14.4:**

$$\text{T}^{bf}(x,z) \leftarrow \text{T}^{bf}(x,y) \wedge \text{T}^{bf}(y,z)$$
$$\Uparrow \qquad \searrow \Uparrow \qquad \searrow$$
$$\text{input}_\text{T}^{bf} \Rightarrow \text{sup}_0[x] \quad \text{sup}_1[x,y] \quad \text{sup}_2[x,z] \Rightarrow \text{output}_\text{T}^{bf}$$

- $\text{sup}_0[x]$ is copied from $\text{input}_\text{T}^{bf}[x]$ (with some exceptions, see exercise)
- $\text{sup}_1[x,y]$ is obtained by joining tables $\text{sup}_0[x]$ and $\text{output}_\text{T}^{bf}[x,y]$
- $\text{sup}_2[x,z]$ is obtained by joining tables $\text{sup}_1[x,y]$ and $\text{output}_\text{T}^{bf}[y,z]$
- $\text{output}_\text{T}^{bf}[x,z]$ is copied from $\text{sup}_2[x,z]$

(we use "named" notation like $[x,y]$ to suggest what to join on; the relations are the same)

## QSQ Evaluation

The set of all auxiliary relations is called a QSQ template (for the given set of adorned rules)

**General evaluation:**

- add new tuples to auxiliary relations until reaching a fixed point
- evaluation of a rule can proceed as sketched on previous slide
- in addition, whenever new tuples are added to a sup relation that feeds into an IDB atom, the input relation of this atom is extended to include all binding given by sup (may trigger subquery evaluation)

$\rightsquigarrow$ there are many strategies for implementing this general scheme

---

**Notation:**

- for an EDB atom $A$, we write $A^{\mathcal{I}}$ for table that consists of all matches for $A$ in the database

---

# Recursive QSQ

Recursive QSQ (QSQR) takes a "depth-first" approach to QSQ

**Evaluation of single rule in QSQR:**

Given: adorned rule $r$ with head predicate $R^\alpha$; current values of all QSQ relations

(1) Copy tuples $\text{input}_R^\alpha$ (that unify with rule head) to $\text{sup}_0^r$

(2) For each body atom $A_1, \ldots, A_n$, do:
   - If $A_i$ is an EDB atom, compute $\text{sup}_i^r$ as projection of $\text{sup}_{i-1}^r \bowtie A_i^{\mathcal{I}}$
   - If $A_i$ is an IDB atom with adorned predicate $S^\beta$:
      (a) Add new bindings from $\text{sup}_{i-1}^r$, combined with constants in $A_i$, to $\text{input}_S^\beta$
      (b) If $\text{input}_S^\beta$ changed, recursively evaluate all rules with head predicate $S^\beta$
      (c) Compute $\text{sup}_i^r$ as projection of $\text{sup}_{i-1}^r \bowtie \text{output}_S^\beta$

(3) Add tuples in $\text{sup}_n^r$ to $\text{output}_R^\alpha$

# QSQR Algorithm

# QSQR Transformation: Example

Predicates S (same generation), p (parent), h (human)

$$S(x, x) \leftarrow h(x)$$
$$S(x, y) \leftarrow p(x, w) \wedge S(v, w) \wedge p(y, v)$$

with query $S(1, x)$.
$\rightsquigarrow$ Query rule: $\text{Query}(x) \leftarrow S(1, x)$

Transformed rules:

$$\text{Query}^f(x) \leftarrow S^{bf}(1, x)$$
$$S^{bf}(x, x) \leftarrow h(x)$$
$$S^{bf}(x, y) \leftarrow p(x, w) \wedge S^{fb}(v, w) \wedge p(y, v)$$
$$S^{fb}(x, x) \leftarrow h(x)$$
$$S^{fb}(x, y) \leftarrow p(x, w) \wedge S^{fb}(v, w) \wedge p(y, v)$$

# Summary and Outlook

Datalog queries can be evaluated bottom-up or top-down

Simplest practical bottom-up technique: semi-naive evaluation

Top-down: Query-Subquery (QSQ) approach (goal-directed)

**Next question:**
- Can bottom-up evaluations be goal directed?
- What about practical implementations?
- Graph databases